

---

# **Introduction to Linear Regression**

*Prof. Nicholas Zabaras  
Center for Informatics and Computational Science*

<https://cics.nd.edu/>

*University of Notre Dame  
Notre Dame, Indiana, USA*

*Email: [nzabaras@gmail.com](mailto:nzabaras@gmail.com)  
URL: <https://www.zabaras.com/>*

*January 28, 2019*



# Contents

---

- [Bayesian Computing and Machine Learning](#), [Motivation to Bayesian inference via a regression example](#), [Over fitting](#), [Effect of Data Size](#), [Training and Test Errors](#), [Over fitting and MLE](#), [Regularization and Model Complexity](#)
- [Prior Modeling](#), [Bayesian Inference and Prediction](#), [Frequentist Vs Bayesian Paradigm](#), [Bias in MLE \(Gaussian Example\)](#)
- [Linear basis function models](#), [MLE and Least Squares](#), [Convexity of the NLL](#), [Sequential Learning of the MLE](#), [Robust Linear Regression](#), [Huber Loss Function](#), [Ridge Regression](#), [Lasso Regularizer and Sparse Solutions](#), [Multi-output Regression](#)
- [The Bias-Variance Decomposition](#)
- [MAP Estimate and Regularized Least Squares](#), [Posterior Distribution](#), [General Gaussian Prior](#), [Sequential MAP Estimation](#), [Example](#), [Predictive Distribution](#), [Pointwise Uncertainty](#), [Plug-in Approximation](#), [Covariance between the Predictions](#), [Basis Functions Vs. Gaussian Processes](#), [Equivalent Kernel Representation](#)
- [Computing the Bias Parameter](#), [Centered Data](#), [Geometry of least squares](#)
  - [Chris Bishop's PRML book](#), Chapters 1 and 2
  - Kevin Murphy's, [Machine Learning: A probabilistic perspective](#), Chapter 5
  - C P Robert, [The Bayesian Choice: From Decision-Theoretic Motivations to Computational Implementation](#), Springer-Verlag, NY, 2001 ([online resource](#))
  - A. Gelman, JB Carlin, HS Stern and DB Rubin, [Bayesian Data Analysis](#), Chapman and Hall CRC Press, 2<sup>nd</sup> Edition, 2003.
  - M Marin and C P Robert, [The Bayesian Core](#), Spring Verlag, 2007 ([online resource](#))
  - Bayesian Statistics for Engineering, [Online Course at Georgia Tech](#), B. Vidakovic.

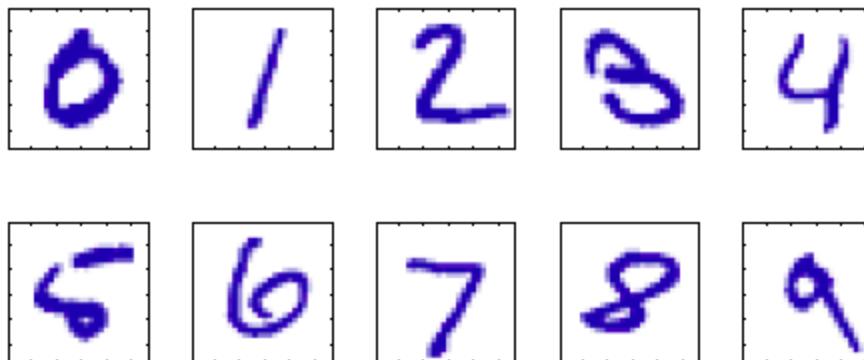


# Typical Problems in Machine Learning

- Pattern Recognition: automatically *classifying the data* into different categories and use of these to take actions.
  - Example: handwritten recognition.

Input: a vector  $x$  of pixel values.

Output: A digit from 0 to 9.



# Typical Problems in Machine Learning

---

- In the digits example, a large set of input vectors  $x_1, \dots, x_N$ , or a *training set* is used to tune the parameters of an adaptive model.
- The category of an input vector is expressed using a target vector  $t$  (identifying the corresponding digit).
- The result of a machine learning algorithm:  $y(x)$  where the output  $y$  is encoded as the target vectors for *any* input  $x$ .



# Terminology

---

- Training or learning phase: determine  $y(x)$  on the basis of the training data.
- Test set, generalization
- Feature extraction
  - Data pre-processing (rotation, scaling, etc.)
  - Using lower-dimensional representation of the input and test data
- Supervised learning (input & target vectors in the training data)
  - Classification (discrete categories) or regression (continuous variables)

# Terminology

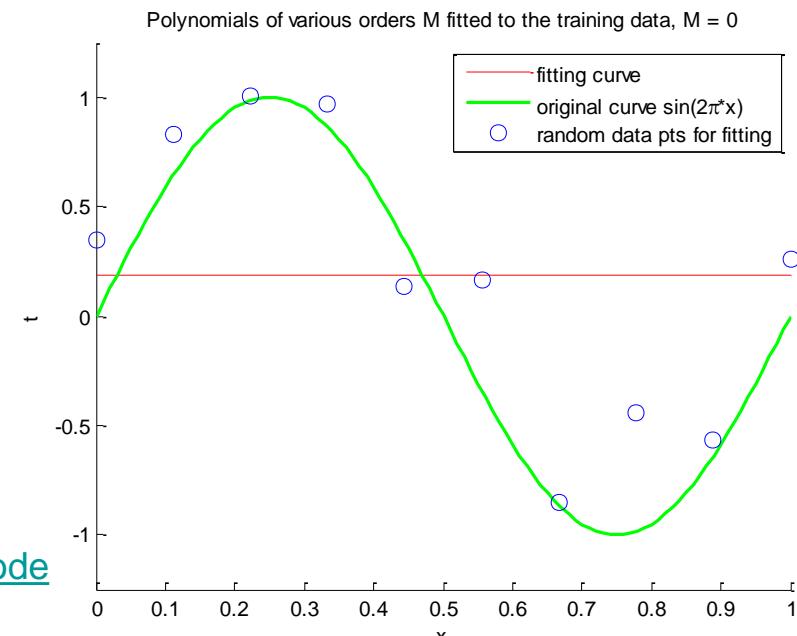
---

- Unsupervised learning (no target vectors in the training data) also called clustering, or density estimation.
- Reinforcement learning ([Richard S. Sutton](#) and [Andrew G. Barto](#))
  - credit assignment (rewards attributed to different moves at the end of a game)
  - exploration (of new actions)
  - exploitation (of high reward actions).

# Motivation Example: Polynomial Curve Fitting

- Problem definition: implicitly trying to discover the underlying (generating) function  $\sin(2\pi x)$  in a set of data.
- Some data points are known:  $x = (x_1, \dots, x_N)^T$  as well as the corresponding target values  $t = (t_1, \dots, t_N)^T$
- We fit the data using a polynomial function of the form

$$\begin{aligned}y(x, w) &= w_0 + w_1 x + \dots + w_M x^M \\&= \sum_{i=0}^M w_i x^i\end{aligned}$$



[MatLab Code](#)



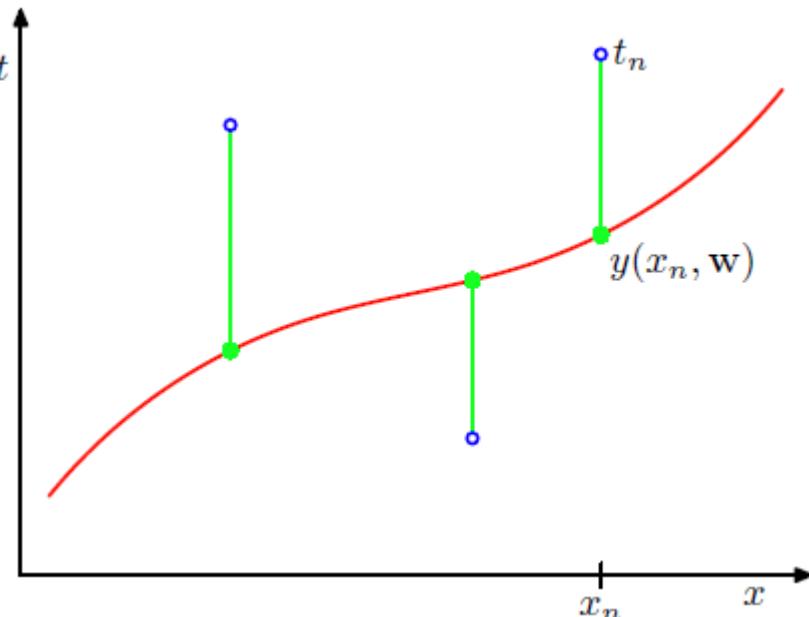
# Motivation Example: Polynomial Curve Fitting

- The values of the coefficients will be determined by fitting the polynomial  $y(x, \mathbf{w}) = \sum_{j=0}^M w_j x^j$  to the training data.
- This can be done by minimizing an error function that measures the misfit between the function  $y(x, \mathbf{w})$ , for any given value of  $\mathbf{w}$ , and the training set data points.

$$\min_{\mathbf{w}} E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2$$

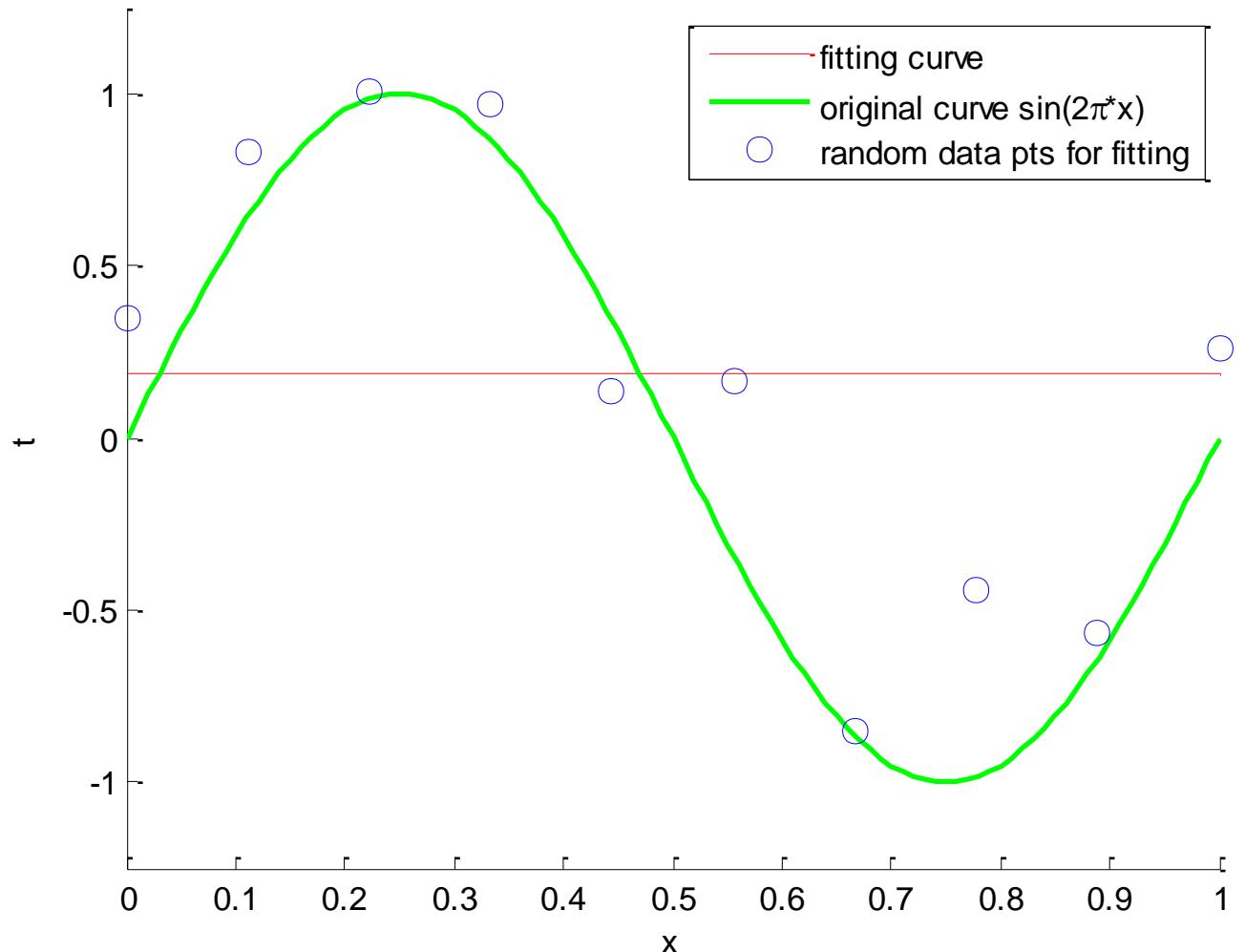
- You can show that the minimizer is obtained by solving:

$$\sum_{j=0}^M A_{ij} w_j = T_i, \quad A_{ij} = \sum_{n=1}^N x_n^{i+j}, \quad T_i = \sum_{n=1}^N x_n^i t_n$$



# Polynomial Curve Fitting

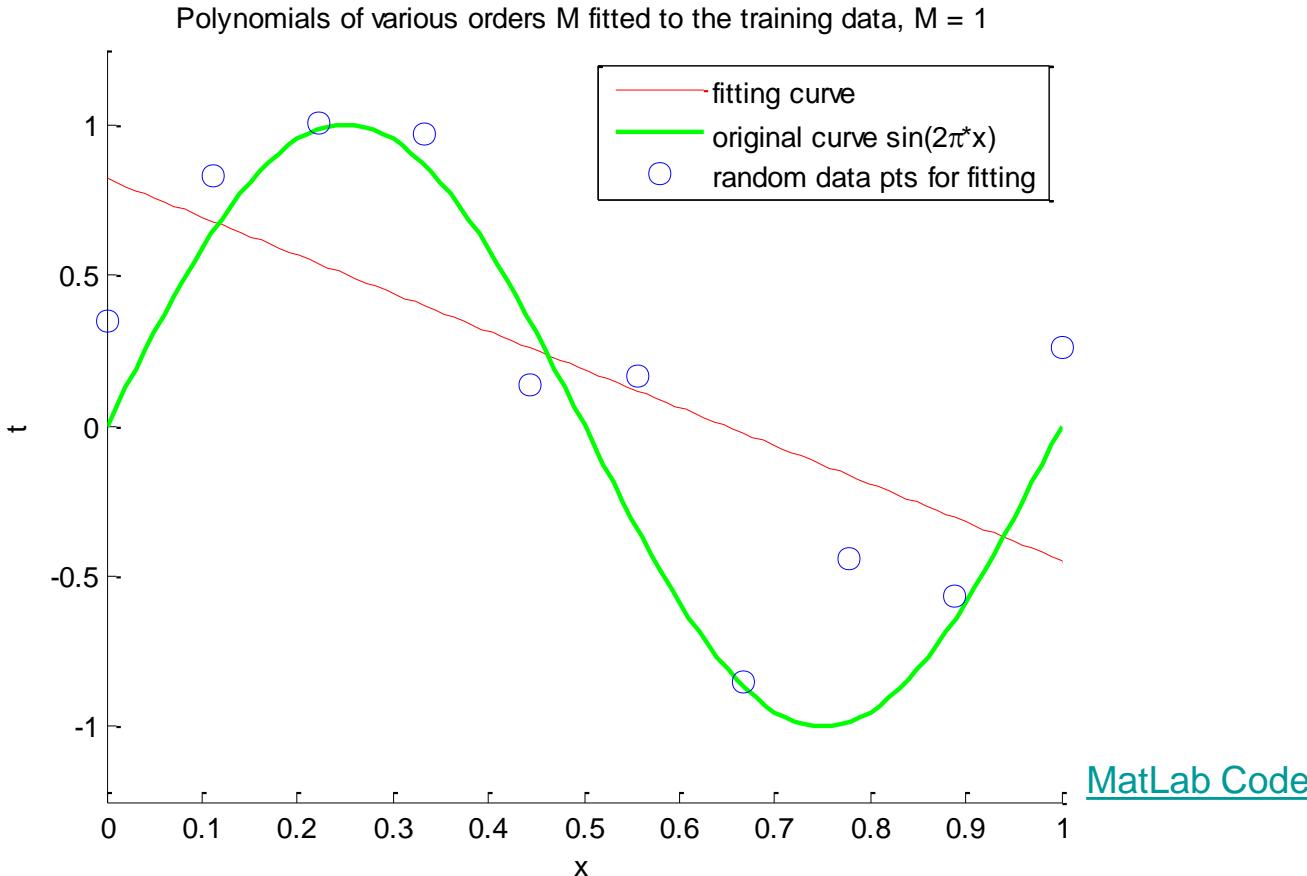
Polynomials of various orders M fitted to the training data, M = 0



[MatLab Code](#)

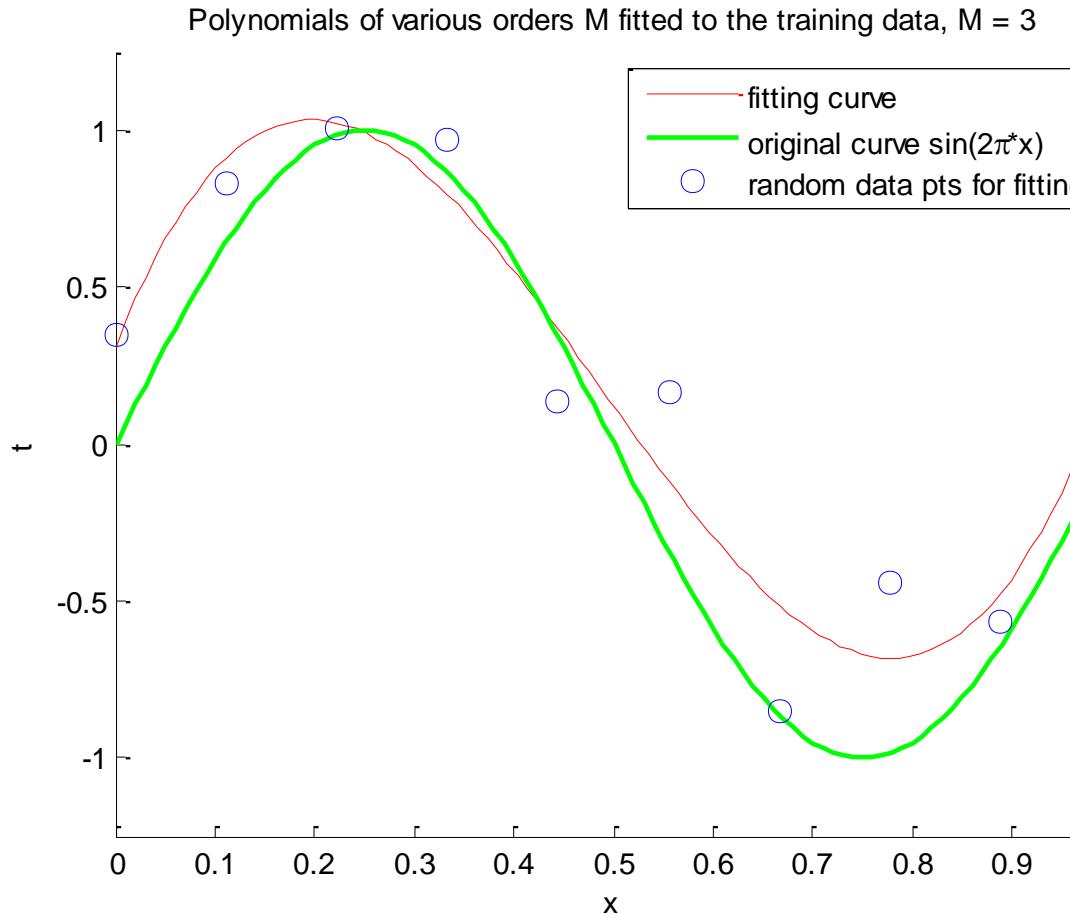


# Polynomial Curve Fitting



- 1<sup>st</sup> order (linear) polynomials give rather poor fit to the data and the  $\sin(2\pi x)$ .

# Polynomial Curve Fitting

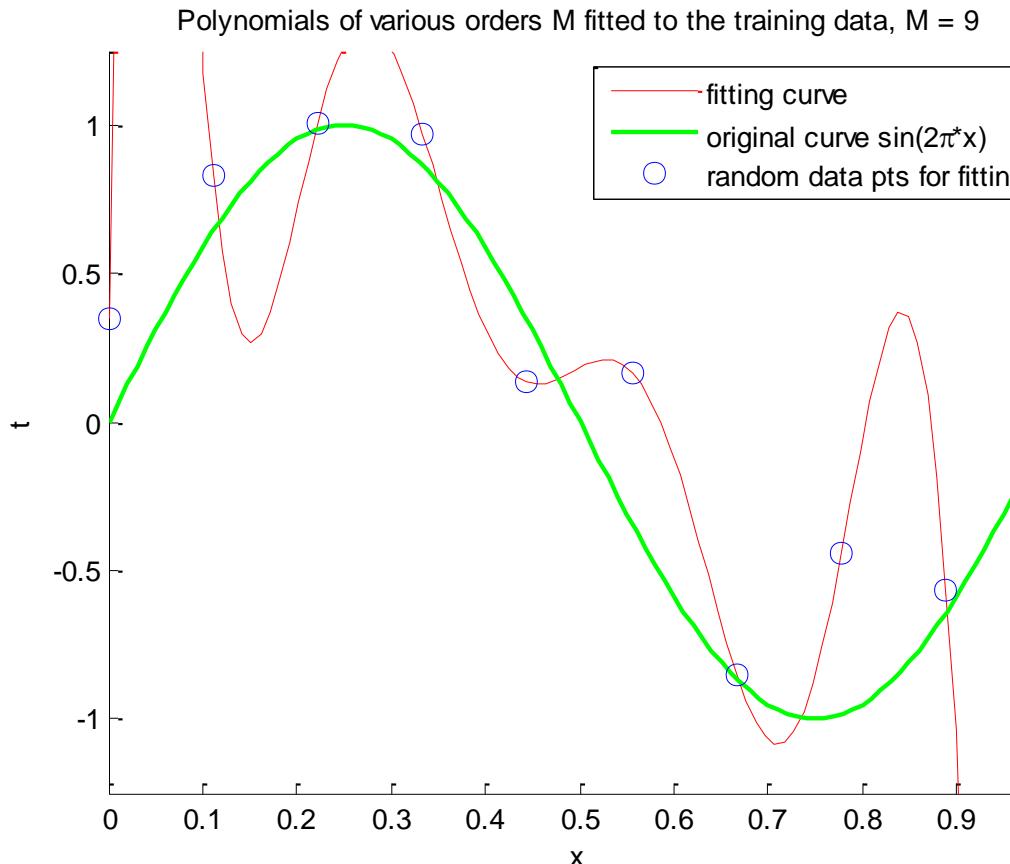


[MatLab Code](#)

- The 3<sup>rd</sup> order polynomial seems to give the best fit to the function  $\sin(2\pi x)$ .



# Overfitting



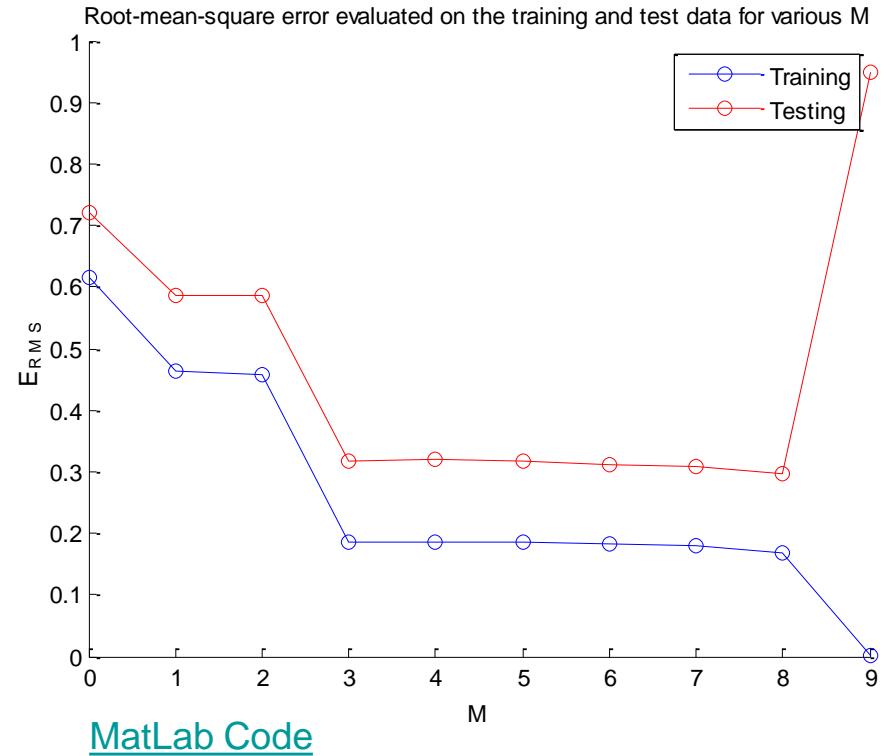
[MatLab Code](#)

- For  $M = 9$  we obtain a perfect fit to the training data. However, the fitted curve oscillates wildly and gives a very poor representation of  $\sin(2\pi x)$ . This is known as **overfitting**.

# Training and Test Errors

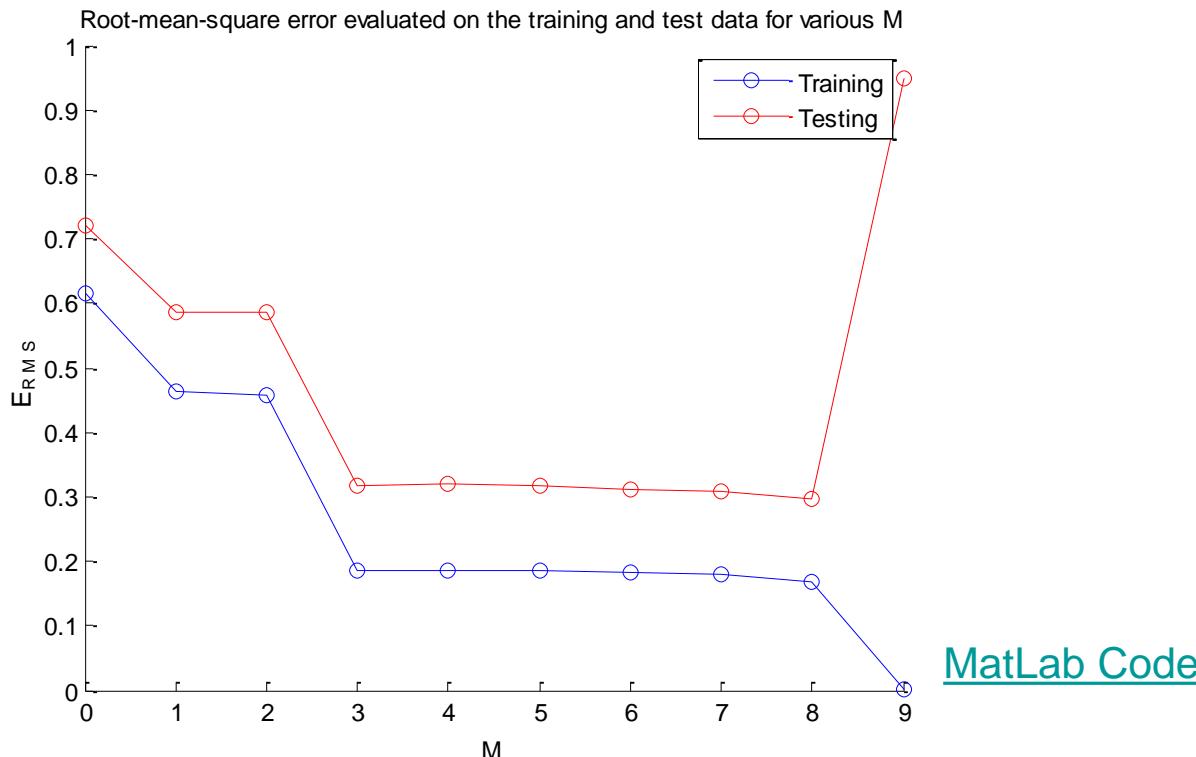
- We often use the root-mean-square (RMS) error. *The division by  $N$  allows us to compare different data sizes. The square root makes  $E_{RMS}$  in units of  $t$ .*

$$E_{RMS} = \left( 2E(w^*) / N \right)^{1/2}$$



# Training and Test Errors

- Small values of  $M$  give relatively large values of the test set error. The corresponding polynomials are incapable of capturing the oscillations in  $\sin(2\pi x)$ .
- $3 < M \leq 8$  give small values for the test set error, and these also give reasonable representation of  $\sin(2\pi x)$ .



# Overfitting

- Obtain insight into the problem by examining the values of  $w$  obtained from polynomials of various order.
- As  $M$  increases, the magnitude of the coefficients typically gets larger.
  - For  $M = 9$ , the coefficients have become finely tuned to the data by developing large positive and negative values.

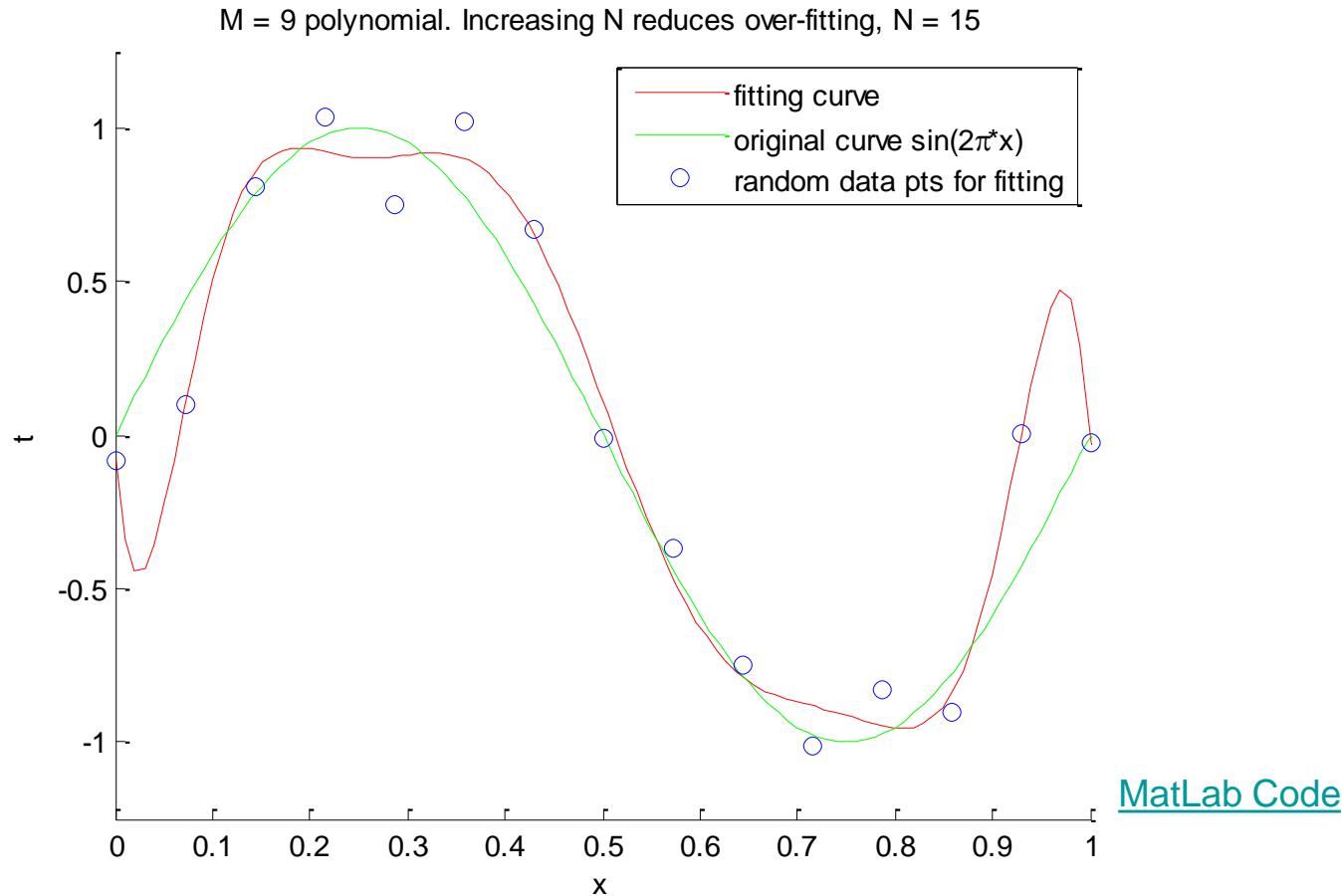
	$M = 0$	$M = 1$	$M = 6$	$M = 9$
$w_0^*$	0.0000	0.0000	0.0000	0.0000
$w_1^*$	0	-0.0000	0.0000	0.0002
$w_2^*$	0	0	-0.0000	-0.0053
$w_3^*$	0	0	0.0000	0.0486
$w_4^*$	0	0	0	-0.2316
$w_5^*$	0	0	0	0.6399
$w_6^*$	0	0	0	-1.0616
$w_7^*$	0	0	0	1.0422
$w_8^*$	0	0	0	-0.5576
$w_9^*$	0	0	0	0.1252

[MatLab Code](#)



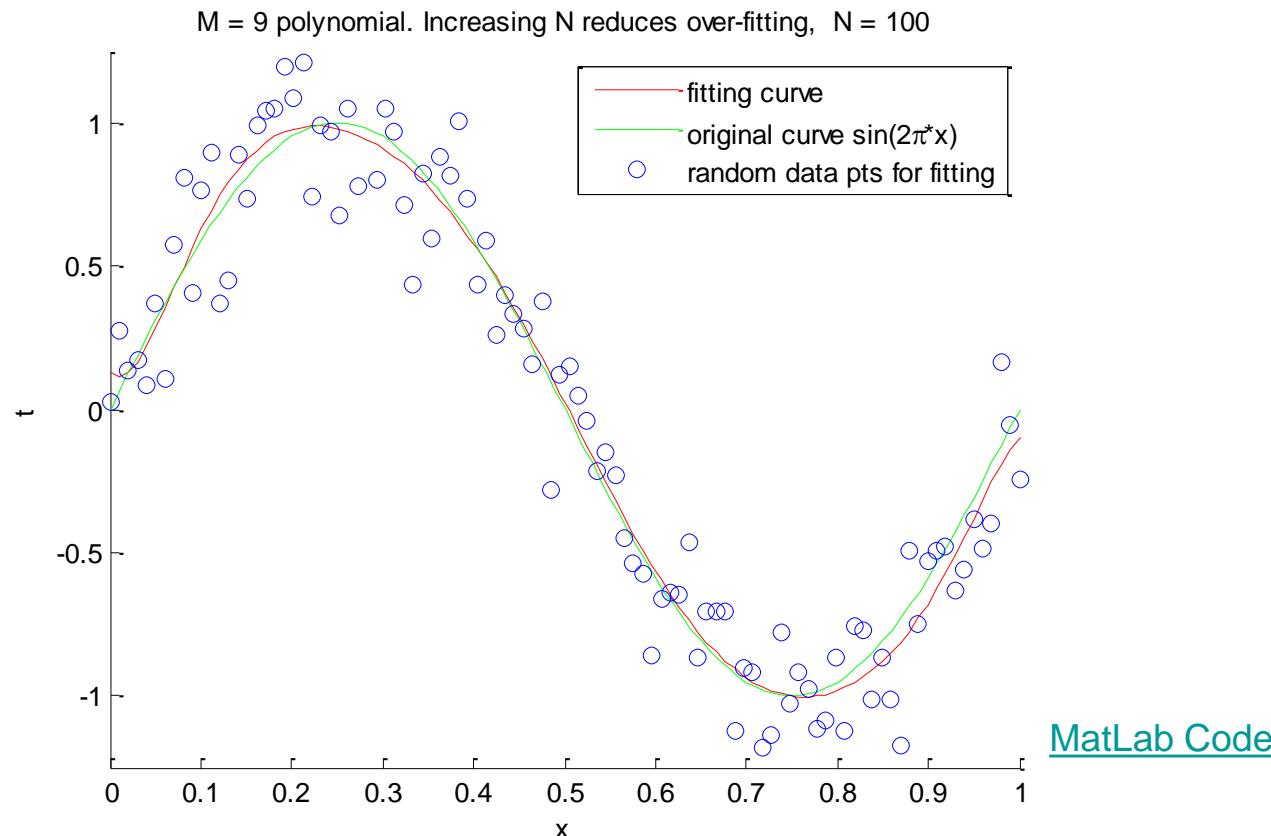
# Varying the Data Size

- It is also interesting to examine the behavior of a given model as the size of the data set is varied.



# Varying the Data Size

- For a given model complexity, the over-fitting problem becomes less severe as the size of the data set increases, i.e. the larger the data set, the more complex (in other words more flexible) the model that we can afford to fit to the data.



# Overfitting and Maximum Likelihood

- It is not reasonable having the number of parameters in a model (model complexity) limited according to the size of the available training set.
- It makes more sense to choose the complexity of the model according to the complexity of the problem being solved.
- The least squares approach to finding the model parameters is a specific case of maximum likelihood.
- The over-fitting problem is a general property of maximum likelihood.

# Overfitting and Bayesian Approach

- By adopting a Bayesian approach, the over-fitting problem can be avoided.
- There is *no difficulty in employing models for which the number of parameters greatly exceeds the number of data points.*
- ***In a Bayesian model, “the effective number” of parameters adapts automatically to the size of the data set.***



# Regularization Technique

- To control the over-fitting we use *regularization*, e.g. adding “a penalty term” to the error function in order to discourage the coefficients from reaching large values.
- The simplest such penalty term takes the form of a sum of squares of all of the coefficients, leading to a modified error function of the form

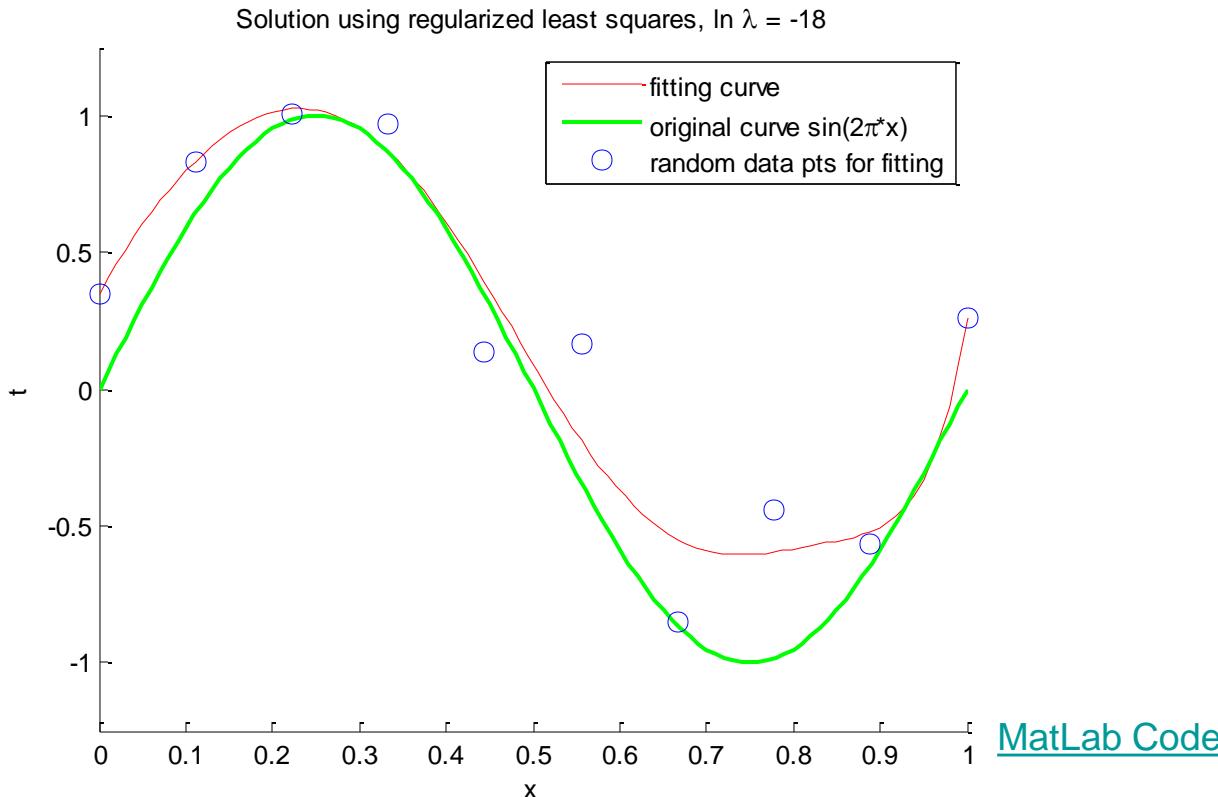
$$\bar{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- The regularization parameter  $\lambda$  governs the relative importance of the regularization term compared with the sum-of-squares error term.
- The minimizer is similar to that given earlier but with

$$\sum_{j=0}^M A_{ij} w_j = T_i, \quad A_{ij} = \sum_{n=1}^N x_n^{i+j} + \lambda \delta_{ij}, \quad T_i = \sum_{n=1}^N x_n^i t_n$$



# Regularization Controls Model Complexity

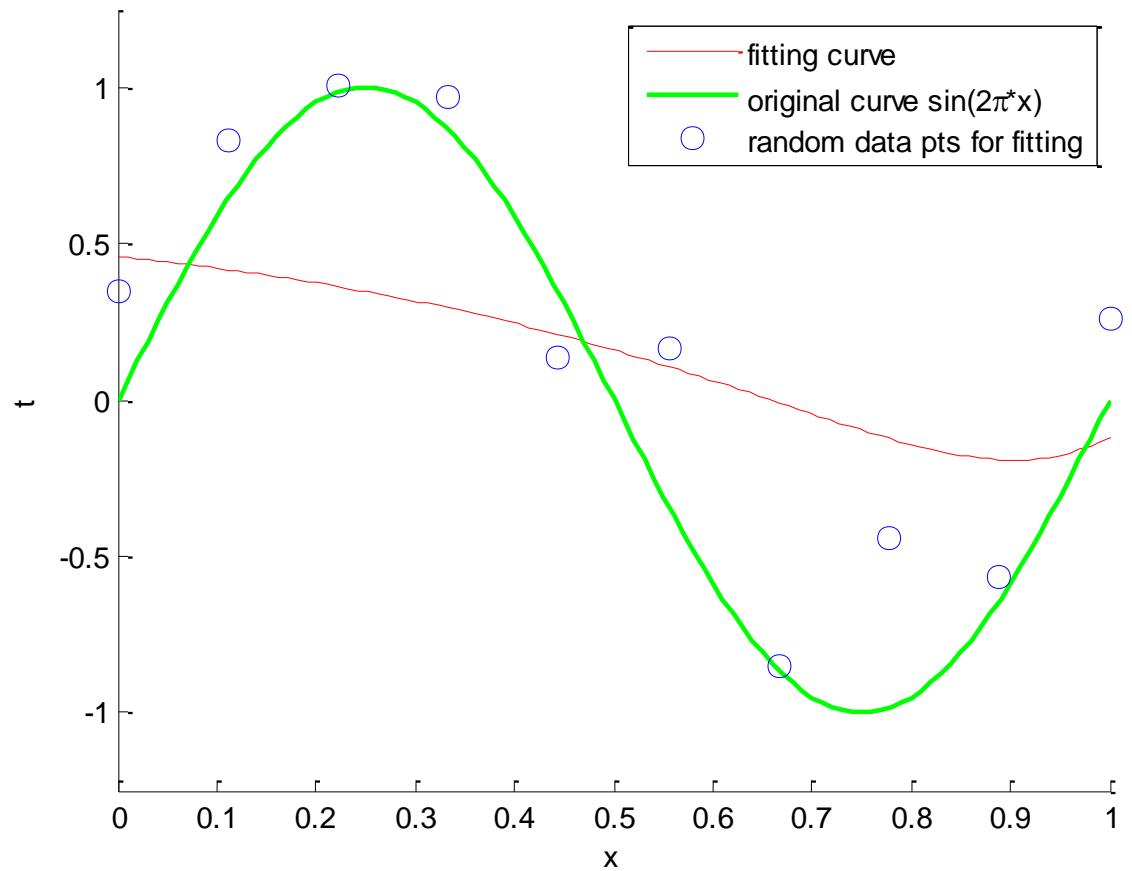


- We now fit the polynomial of order  $M = 9$  to the same data set as before but now using the regularized error function.
- For  $\ln l = -18$ , the over-fitting has been suppressed and we obtain a much closer representation of  $\sin(2\pi x)$ .



# Regularization Controls Model Complexity

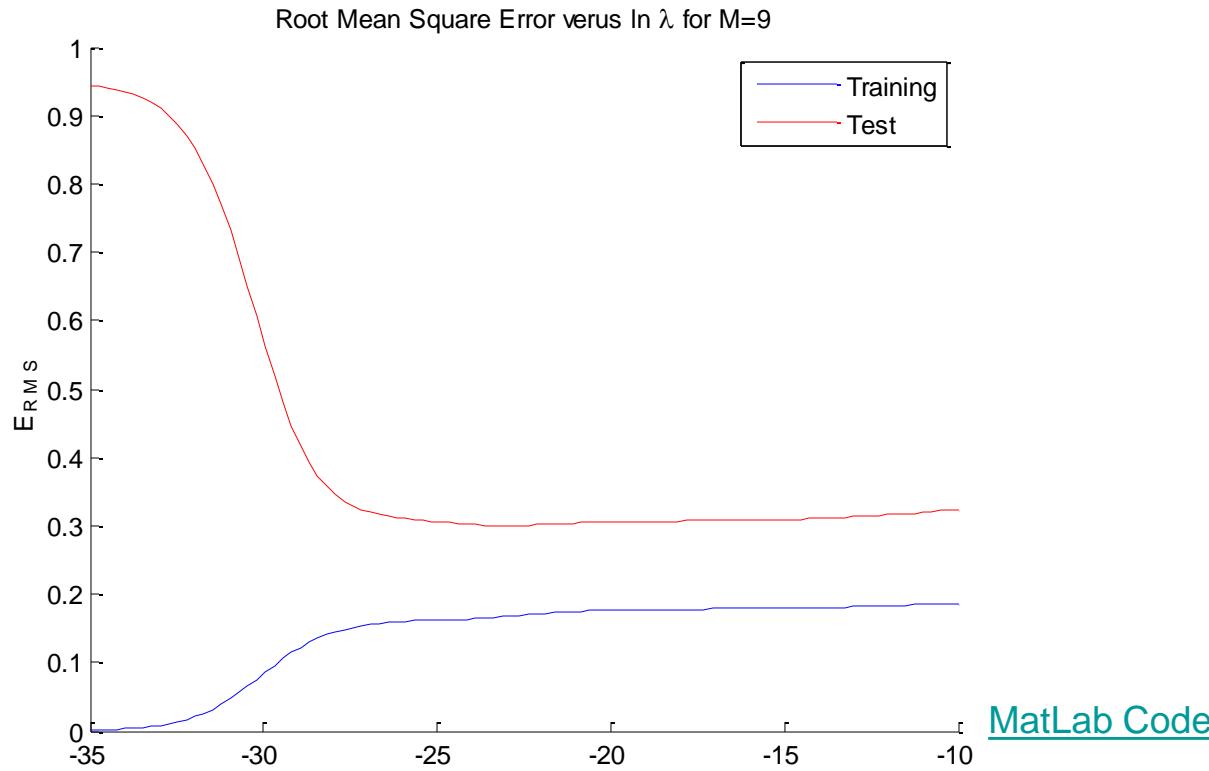
Solution using regularized least squared,  $\ln \lambda = 0$



[MatLab Code](#)

- If, however, we use too large a value for  $\lambda$  then we again obtain a poor fit, as shown for  $\ln\lambda = 0$ .

# Regularization Controls Model Complexity



- ❑  $\lambda$  controls the effective complexity of the model and hence determines the degree of over-fitting.
- ❑ We will soon re-examine this problem with a Bayesian approach that avoids the over-fitting problem.

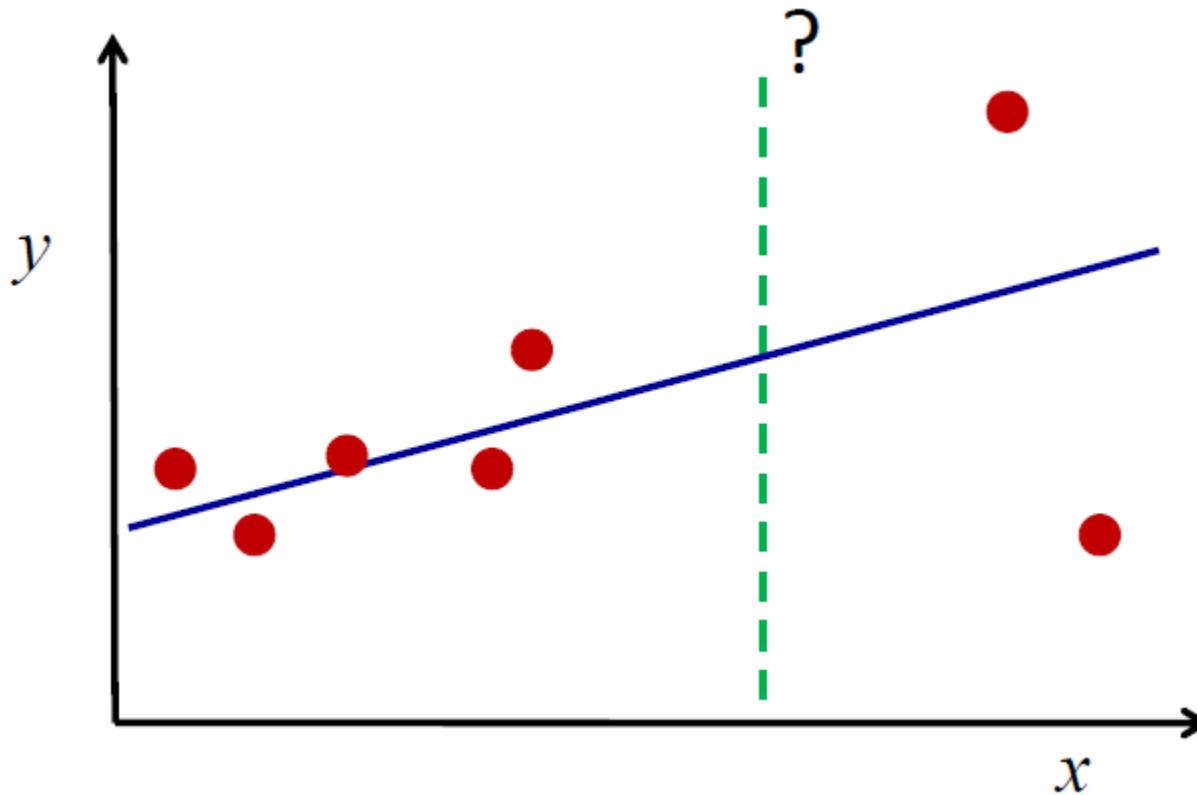


# MLE, Regularization and Model Complexity

- Effective model complexity in MLE is governed by the number of basis functions and is controlled by the size of the data set.
- With regularization, the effective model complexity is controlled mainly by  $\lambda$  and still by the number and form of the basis functions.
- *Thus the model complexity for a particular problem cannot be decided simply by maximizing the likelihood function as this leads to excessively complex models and over-fitting.*
- Independent hold-out data can be used to determine model complexity but this wastes data and it is computationally expensive.

# Prior Knowledge is Essential

- We cannot do everything simply based on data – prior knowledge is essential to inference and prediction.



# Frequentist Versus Bayesian Paradigms

- The likelihood  $p(\mathcal{D}|w)$  is essential in both Bayesian and frequentist approaches but it is used in different roles.
- In a frequentist approach
  - $w$  is a fixed parameter computed by an estimator (e.g. maximum likelihood estimator).
  - Error bars on this point estimate are computed by considering the distribution of all possible data sets  $\mathcal{D}$  (e.g. variability of predictions between different bootstrap data sets)
- In Bayesian approach
  - there is only one set of data  $\mathcal{D}$  and
  - the uncertainty in  $w$  is introduced with appropriate prior and computing posterior probabilities over  $w$ .



# Bayesian Probabilities

---

- For example in the regression problem with the observed data  $\mathcal{D} = \{t_1, \dots, t_N\}$ , we can obtain the conditional probability  $p(\mathbf{w}|\mathcal{D})$  by Bayes' theorem

$$p(\mathbf{w} | \mathcal{D}) = \frac{p(\mathcal{D} | \mathbf{w}) p(\mathbf{w})}{p(\mathcal{D})}, \quad p(\mathcal{D}) = \int p(\mathcal{D} | \mathbf{w}) p(\mathbf{w}) d\mathbf{w}$$

- The quantity  $p(\mathcal{D}|\mathbf{w})$  on the right-hand side of Bayes' theorem is evaluated for the observed data set  $\mathcal{D}$  and can be viewed as a function of the parameter vector  $\mathbf{w}$  (*likelihood function*).
- Given this definition of likelihood, we can state Bayes' theorem in words

$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$



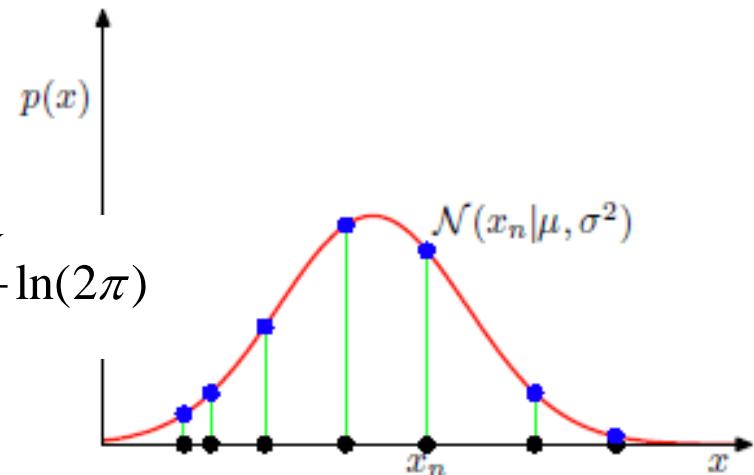
# Gaussian Distribution

- Consider the Gaussian distribution

$$\mathcal{N}(x | \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\}$$

- The likelihood function for the Gaussian distribution is

$$p(\mathcal{D} | \mu, \sigma^2) = \prod_{n=1}^N \mathcal{N}(x_n | \mu, \sigma^2)$$



- The log likelihood takes the form\*

$$\ln p(\mathcal{D} | \mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi)$$

- Maximum likelihood solution

$$\mu_{ML} = \frac{\sum_{n=1}^N x_n}{N}, \quad \sigma_{ML}^2 = \frac{\sum_{n=1}^N (x_n - \mu_{ML})^2}{N}$$

\* We work often with log-likelihood to avoid underflow (taking products of small probabilities) and for simplifying the algebra.

# MLE for a Gaussian Distribution

$$\mu_{ML} = \underbrace{\frac{1}{N} \sum_{i=1}^N x_i}_{\text{Sample mean}}, \sigma_{ML}^2 = \underbrace{\frac{1}{N} \sum_{i=1}^N (x_i - \mu_{ML})^2}_{\text{Sample variance wrt ML mean (not the exact mean)}}$$

- The MLE approach underestimates the variance (bias) – this is at the root of the over-fitting problem e.g. in polynomial curve fitting.
- The maximum likelihood solutions  $\mu_{ML}, \sigma_{ML}^2$  are functions of the data set values  $x_1, \dots, x_N$ . Consider the expectations of these quantities with respect to the data set values, which come from a Gaussian.
- Using the point estimates above you can show that :

$$\mathbb{E}[\mu_{ML}] = \mu, \quad \mathbb{E}[\sigma_{ML}^2] = \frac{N-1}{N} \sigma^2$$

*In this derivation*

*you need to use :*

$$\mathbb{E}[x_i x_j] = \mu^2 \text{ for } i \neq j$$

$$\mathbb{E}[x_i^2] = \sigma^2 + \mu^2$$

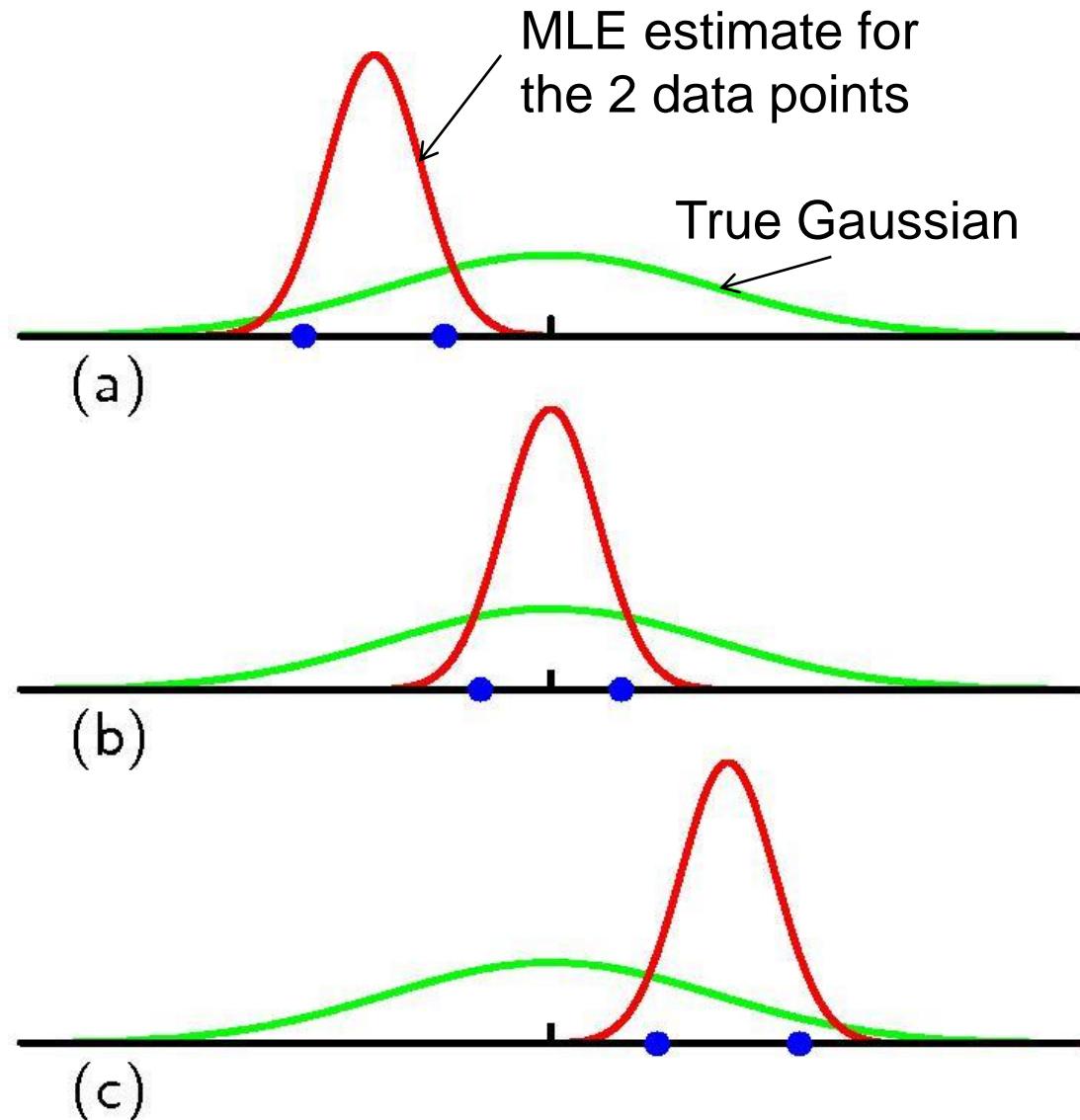


# MLE: Underestimating the Variance

- In the schematic from [Bishop's PRML](#), we consider 3 cases each with 2 data points extracted from the true Gaussian.

- The mean of the 3 distributions predicted via MLE (i.e. averaged over the data) is correct.

- However, ***the variance is underestimated*** since it is a variance with respect to the sample mean and NOT the true mean.



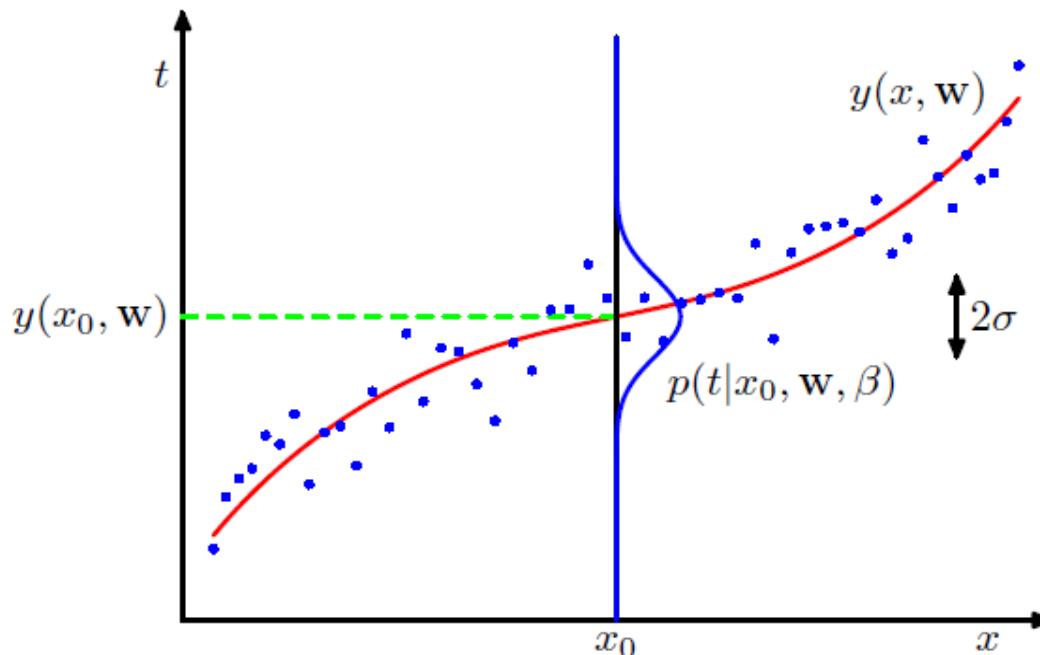
# Linear Basis Functions Models

---

- We are now interested in a linear combination – regression – of a ``fixed set'' of nonlinear basis functions.
- Supervised learning:  $N$  observations  $\{\mathbf{x}_n\}$  with corresponding target values  $\{t_n\}$  are provided.
- The goal is to predict  $t$  for a new value  $\mathbf{x}$ .
- We construct a function such that  $y(\mathbf{x})$  is a prediction of  $t$ .
- We follow a Bayesian perspective and model the predictive distribution  $p(t|\mathbf{x})$ .

# Linear Regression

- From the conditional distribution  $p(t|x)$ , we can make point estimates of  $t$  for a given  $x$  by minimizing a ‘loss function’.
- For a quadratic loss function, the point estimate is the conditional mean  $y(x, w) = \mathbb{E}[t|x]$ .



# Linear Regression

---

- The simplest linear model for regression is one that involves a linear combination of the input variables

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_D x_D$$

where  $\mathbf{x} = (x_1 \quad x_2 \quad \dots \quad x_D)^T$ .

- This is often simply known as linear regression.
- $D$  is the input dimensionality.
- This model is a linear function of the parameters.



# Linear Basis Functions Models

- More generally regression function is:

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1\phi_1(\mathbf{x}) + \dots + w_{M-1}\phi_{M-1}(\mathbf{x}) = \sum_{i=0}^{M-1} w_i\phi_i(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

where  $\phi_i(\mathbf{x})$  are known as basis functions and

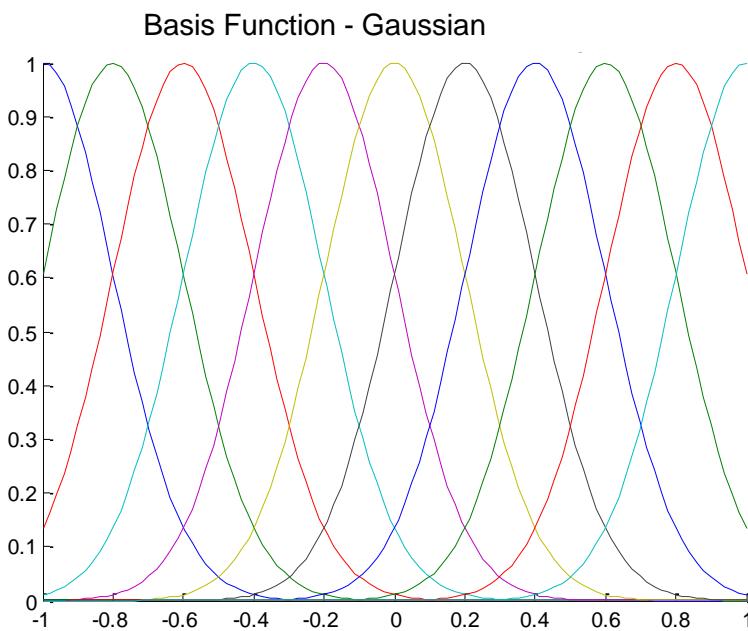
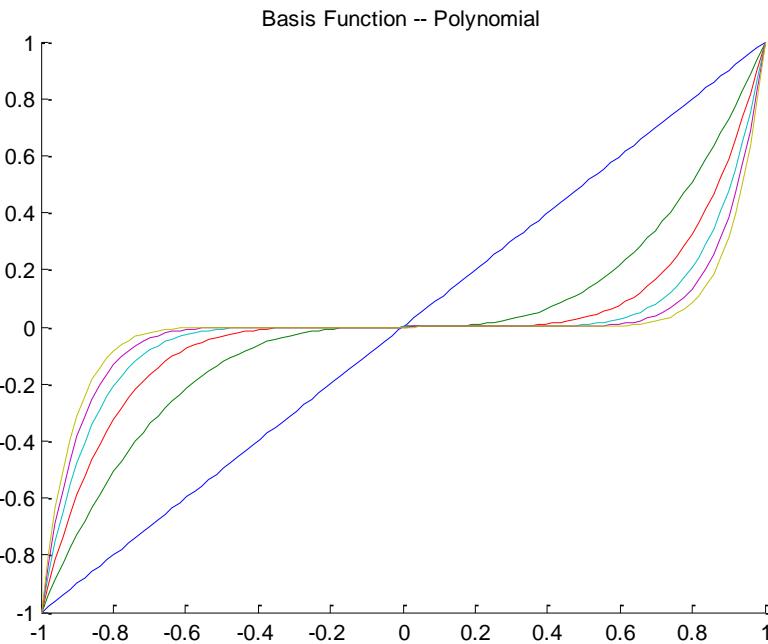
$$\boldsymbol{\phi}(\mathbf{x}) = (\phi_0(\mathbf{x}) \quad \phi_1(\mathbf{x}) \quad \phi_2(\mathbf{x}) \quad \dots \quad \phi_{M-1}(\mathbf{x}))^T, \phi_0(\mathbf{x}) = 1, w_0 = bias$$

- Often  $\{\phi_i(\mathbf{x})\}$  represent features extracted from the data .

$\mathbf{x}$



# Polynomial and Gaussian Basis Functions



[MatLab code](#)

$$\mu = -1 : 0.2 : 1$$

$$s = 0.2$$

Gaussian basis functions:

$$\phi_j(x) = \exp\left(-\frac{(x - \mu_j)^2}{2s^2}\right)$$

Polynomial basis functions  
(scalar input, global support):

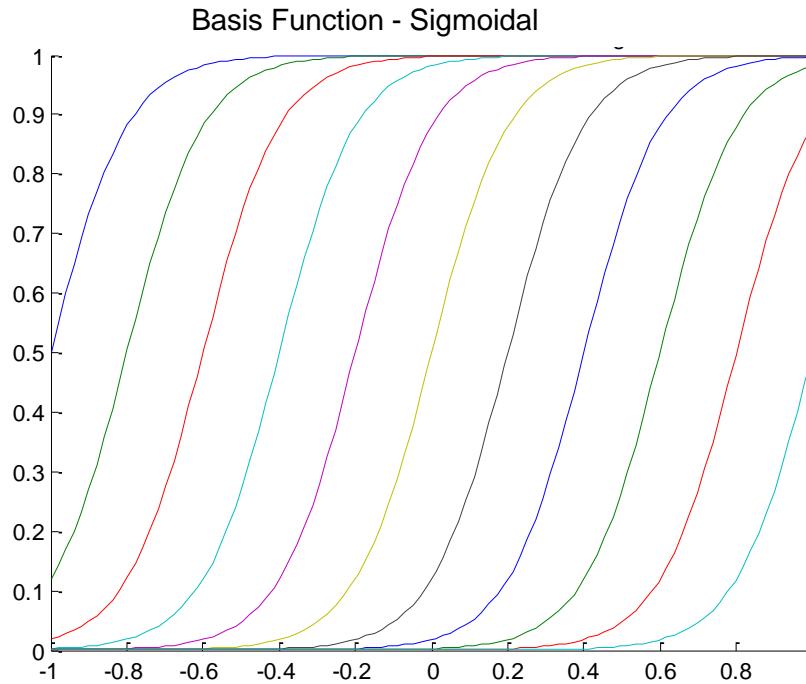
$$\phi_j(x) = x^j$$



# Logistic Sigmoidal Basis Functions

$$\mu = -1 : 0.2 : 1$$

$$s = 0.1$$



[MatLab code](#)

Sigmoidal basis functions:

$$\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right), \quad \sigma(a) = \frac{1}{1 + e^{-a}}$$

$$\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}} = 2\sigma(2a) - 1$$

$\sigma(a)$ : logistic sigmoid function

# Sigmoidal and Tanh Basis Functions

- The sigmoidal and tanh basis functions are related:

$$\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}} = 2\sigma(2a) - 1 \quad \text{where} \quad \sigma(a) = \frac{1}{1 + e^{-a}}$$

- A general linear combination of logistic sigmoidal functions is equivalent to a linear combination of tanh functions:

$$\begin{aligned} y(x, w) &= w_0 + \sum_{j=1}^{M-1} w_j \sigma\left(\frac{x - \mu_j}{s}\right) = w_0 + \sum_{j=1}^{M-1} w_j \sigma\left(2 \frac{x - \mu_j}{2s}\right) \\ &= w_0 + \sum_{j=1}^{M-1} w_j \frac{1 + \tanh\left(\frac{x - \mu_j}{2s}\right)}{2} = u_0 + \sum_{j=1}^{M-1} u_j \tanh\left(\frac{x - \mu_j}{2s}\right) \end{aligned}$$

$$\text{where: } u_0 = w_0 + \frac{1}{2} \sum_{j=1}^{M-1} w_j, \quad u_j = \frac{w_j}{2}$$



# Choice of Basis Functions

---

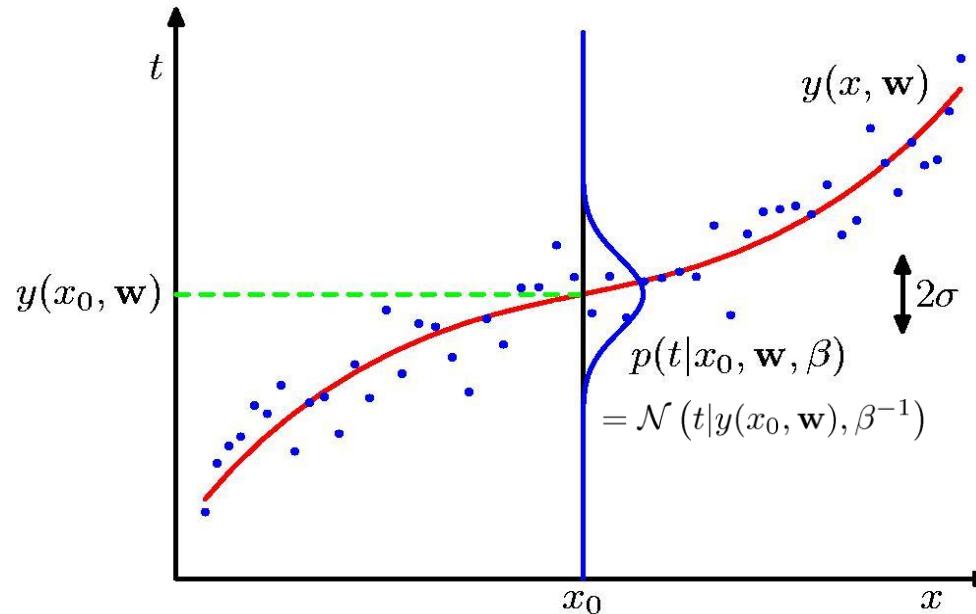
- We are interested in functions of *local support* to explore adaptivity.
- Local support functions comprise a spectrum of different spatial frequencies.
- An example is wavelets that are local both spatially and in frequency.
  - They are however useful only when the input is defined on a lattice.

# MLE and Least Squares

- Consider a set of training data comprising  $N$  input  $x = (x_1, \dots, x_N)^T$  & the corresponding target values  $t = (t_1, \dots, t_N)^T$
- We assume that, given the value of  $x$ , the corresponding value of  $t$  has a Gaussian distribution with a mean equal to the value  $y(x, w)$  and precision (inverse of the variance)  $\beta$

$$p(t | x, w, \beta) = \mathcal{N}(t | y(x, w), \beta^{-1})$$

$$\begin{aligned}y(x, w) &= \phi(x)^T w \\t &= y(x, w) + \varepsilon \\ \varepsilon &\sim \mathcal{N}(0, \beta^{-1})\end{aligned}$$



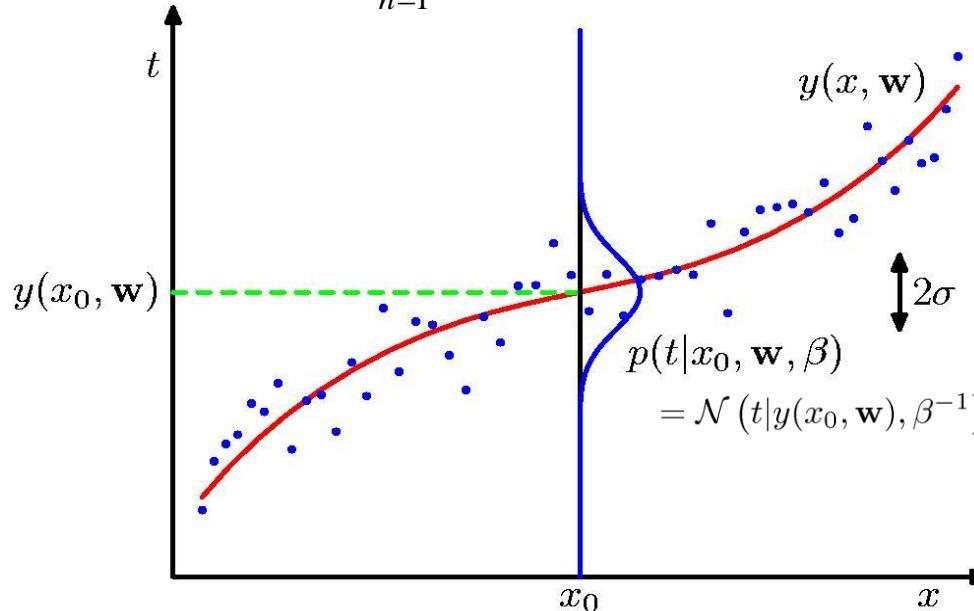
# MLE and Least Squares

- The likelihood function is

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y(x_n, \mathbf{w}), \beta^{-1})$$

- From this, the log-likelihood takes the form:

$$\ln p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$



# **MLE and Least Squares**

---

- Consider first the MLE estimate for  $\mathbf{w}$ . Note that maximizing the log-likelihood to obtain  $\mathbf{w}_{ML}$  is the same as minimizing the sum of squares error function (residual sum of squares,  $RSS(\mathbf{w})$ ):

$$\max_{\mathbf{w}} \ln p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) \Leftrightarrow$$

$$\mathbf{w}_{ML} = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

- We can also determine the MLE estimate of  $\beta$ :

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, \mathbf{w}_{ML}) - t_n\}^2$$

- We can now make (a frequentist, plug-in approximation) prediction as follows:  $p(t | \mathbf{x}, \mathbf{w}_{ML}, \beta_{ML}) = \mathcal{N}(t | y(\mathbf{x}, \mathbf{w}_{ML}), \beta_{ML}^{-1})$

# Maximum Likelihood and Least Squares

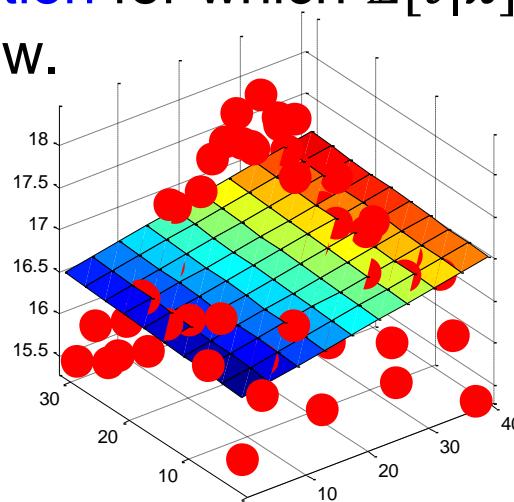
- Assume observations from a deterministic function with added Gaussian noise:

$$t = y(\mathbf{x}, \mathbf{w}) + \varepsilon \quad \text{where} \quad p(\varepsilon | \beta) = \mathcal{N}(\varepsilon | 0, \beta^{-1})$$

which is the same as saying,

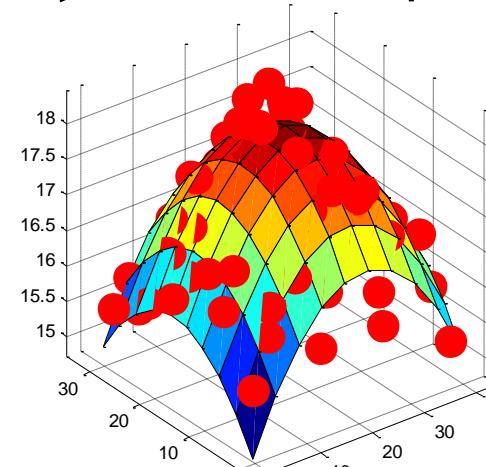
$$p(t | \mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t | y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

- Here  $\beta$  is the precision. This is based on a squared loss function for which  $\mathbb{E}[t|\mathbf{x}] = y(\mathbf{x}, \mathbf{w})$ . A 2D example is shown below.



$$\mathbb{E}[t | \mathbf{x}] = w_0 + w_1x_1 + w_2x_2$$

Machine Learning, University of Notre Dame, Notre Dame, IN, USA (Spring 2019, N. Zabaras)



$$\mathbb{E}[t | \mathbf{x}] = w_0 + w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2$$

Run [surfaceFitDemo](#) from PMTK3

# Maximum Likelihood and Least Squares

- Given observed inputs,  $X = \{x_1, \dots, x_N\}$ , and targets  $t = \{t_1, \dots, t_N\}^T$ , we obtain the likelihood function

$$p(t | X, w, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | w^T \phi(x_n), \beta^{-1})$$

- We often use the *log-likelihood*:

$$\ell(w, \beta) = \log p(\mathcal{D} | \theta) = \sum_{n=1}^N \log \mathcal{N}(t_n | w^T \phi(x_n), \beta^{-1}), \theta = (w, \beta)$$

- Instead of maximizing the log-likelihood, one equivalently can minimize *the negative log-likelihood (NLL)*

$$NLL(w, \beta) = -\sum_{n=1}^N \log \mathcal{N}(t_n | w^T \phi(x_n), \beta^{-1})$$

# Maximum Likelihood and Least Squares

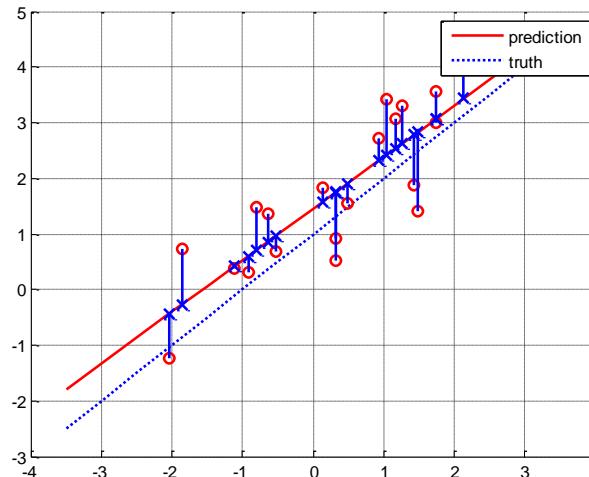
- Taking the log of the likelihood, we obtain:

$$\ln p(t | w, \beta) = \sum_{n=1}^N \ln \mathcal{N}(t_n | w^T \phi(x_n), \beta^{-1}) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(w)$$

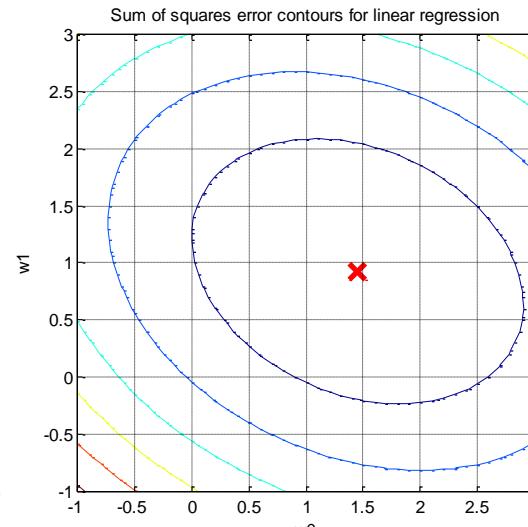
where with  $E_D(w)$  we have denoted:

$$E_D(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w^T \phi(x_n))^2, \text{ RSS} = \sum_{n=1}^N (t_n - w^T \phi(x_n))^2, \text{ MSE} = \text{RSS} / N$$

- RSS is often known as the *residual sum of squares* or *sum of squared errors (SSE)* and *MSE* is the mean squared error.
- Computing  $w$  via *MLE* is the same as *Least Squares*.



Run [residualsDemo](#) from PMTK3



The NLL is a quadratic bowl with a unique minimum (the MLE estimate)

Run [contoursSSEdemo](#), from PMTK3



# Maximum Likelihood and Least Squares

$$\ln p(\mathbf{t} | \mathbf{w}, \beta) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w}), E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2$$

- Setting the gradient of the log-likelihood (written here as a row vector) wrt  $\mathbf{w}$  equal to zero:

$$\nabla_{\mathbf{w}} \ln p(\mathbf{t} | \mathbf{w}, \beta) = \beta \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n)^T = \beta \left\{ \sum_{n=1}^N t_n \phi(\mathbf{x}_n)^T - \mathbf{w}^T \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \right\} = 0$$

- This equation can be solved for  $\mathbf{w}$ .



# Maximum Likelihood and Least Squares

$$\mathbf{w}^T \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T = \sum_{n=1}^N t_n \phi(\mathbf{x}_n)^T \Rightarrow \mathbf{w}^T \Phi^T \Phi = (\Phi^T \mathbf{t})^T \Rightarrow \Phi^T \Phi \mathbf{w} = \Phi^T \mathbf{t}$$

□ We obtain (normal equation; ordinary least squares solution):

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \equiv \Phi^\dagger \mathbf{t}, \quad \Phi^\dagger : \text{Moore-Penrose pseudo-inverse}$$

where we have defined:

$$\Phi = \begin{pmatrix} \phi^T(\mathbf{x}_1) \\ \phi^T(\mathbf{x}_2) \\ \vdots \\ \phi^T(\mathbf{x}_N) \end{pmatrix} = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & .. & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & .. & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \vdots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & .. & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}, \quad \begin{aligned} \Phi^T &= (\phi(\mathbf{x}_1) \quad \phi(\mathbf{x}_2) \quad .. \quad \phi(\mathbf{x}_N)) \\ \phi(\mathbf{x}_i) &\equiv (\phi_0(\mathbf{x}_i) \quad \phi_1(\mathbf{x}_i) \quad .. \quad \phi_{M-1}(\mathbf{x}_i))^T \\ \Phi &= (\varphi_0 \quad \varphi_1 \quad .. \quad \varphi_{M-1}) \\ \varphi_i &\equiv (\phi_i(\mathbf{x}_1) \quad \phi_i(\mathbf{x}_2) \quad .. \quad \phi_i(\mathbf{x}_N))^T \end{aligned}$$

□ Note that indeed:

$$\Phi^T \Phi = \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T, \quad \Phi^T \mathbf{t} = \sum_{n=1}^N t_n \phi(\mathbf{x}_n)$$



# Maximum Likelihood and Least Squares

$$\ln p(\mathbf{t} | \mathbf{w}, \beta) = \sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w})$$

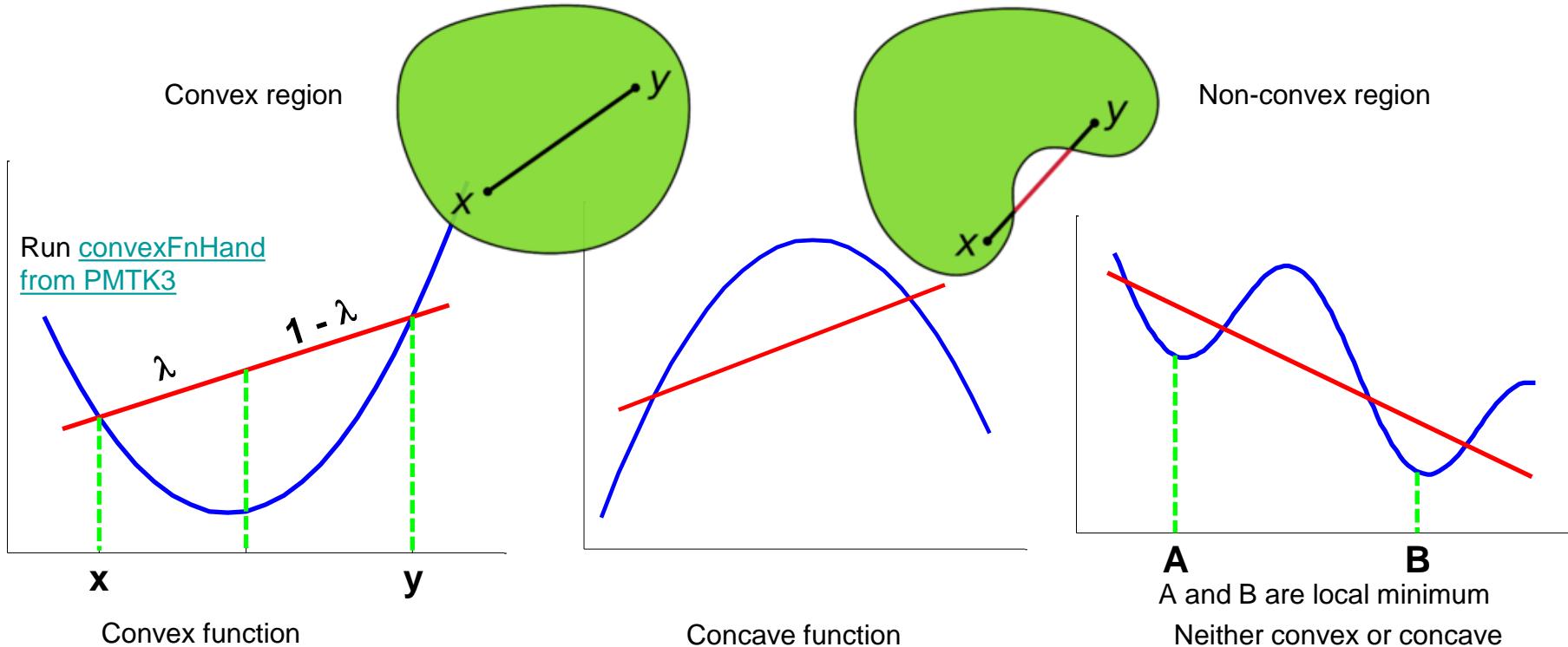
- Taking now the log of the likelihood with respect to  $\beta$  gives:

$$\frac{1}{\beta_{ML}} = \frac{2}{N} E_D(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2$$

- So the MLE variance  $\beta_{ML}$  is equal to the residual variance of the target values around the regression function.

# Convexity of the NLL

- Convexity of the NLL (positive definite Hessian) leads to a unique globally optimal MLE.
- Some models of interest don't have concave likelihoods and locally optimal MLE estimates are needed.



$$\theta^2, e^\theta, \theta \log \theta (\theta > 0)$$

Machine Learning, University of Notre Dame, Notre Dame, IN, USA (Spring 2019, N. Zabaras)

# Sequential Learning: LMS Algorithm

- If the data set is large, we use sequential (on-line) algorithms
- We apply the technique of ***stochastic (sequential) gradient descent.***
- If the error function comprises a sum over data points  $E = \sum_n E_n$

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2$$

then after presentation of pattern  $n$ , the stochastic gradient descent algorithm updates the parameter vector  $\mathbf{w}$  using

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n = \mathbf{w}^{(\tau)} + \underbrace{\eta \left( t_n - \mathbf{w}^{(\tau)^T} \phi(\mathbf{x}_n) \right) \phi(\mathbf{x}_n)}_{-\nabla E_n}$$

$\tau$  is the iteration number &  $\eta$  the learning rate parameter.

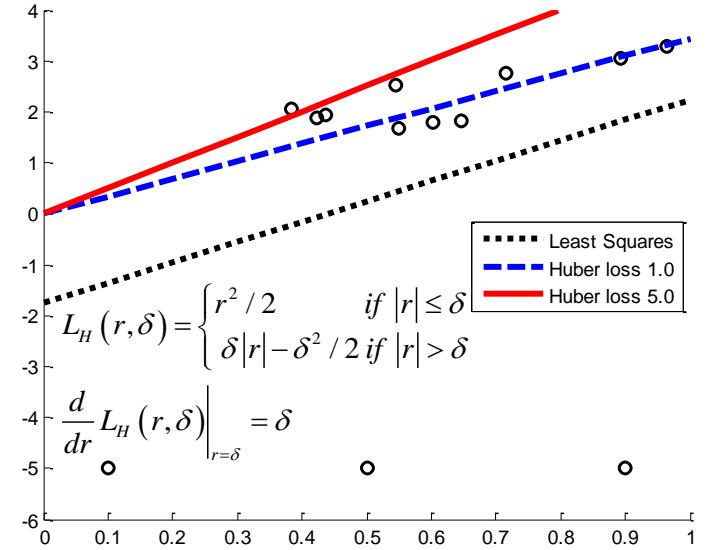
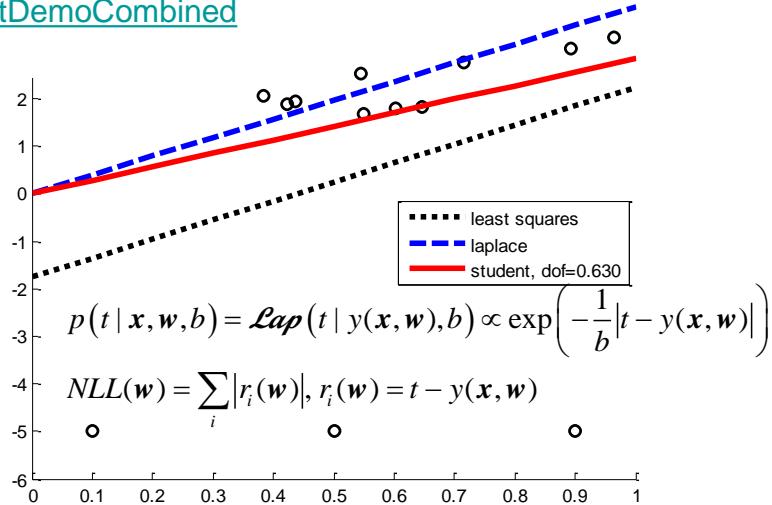
- This is known as ***least-mean-squares*** or the ***LMS algorithm.***



# Robust Linear Regression

- Using a Gaussian distribution for the noise,  
$$t = y(\mathbf{x}, \mathbf{w}) + \varepsilon, \varepsilon \sim \mathcal{N}(\varepsilon | 0, \beta^{-1})$$
 can result in poor fit especially if we have outliers in the data.
- Squared error penalizes deviations quadratically, so points far from the line have more affect on the fit than points near the line.
- To achieve robustness to outliers one can replace the Gaussian with a distribution that has heavy tails (e.g. the Laplace distribution). Such a distribution assigns higher likelihood to outliers, without having to perturb the regression line to “explain” them.

Run [linregRobustDemoCombined](#)  
from PMTK3



# Robust Linear Regression

- Using the Laplace distribution leads to  $L_1$  error norm (non-linear objective function) that is difficult to optimize.
- A solution is to transform the problem (by increasing its dimension to  $2N + M$ ) to a *linear programming problem*.

$$r_i \triangleq r_i^+ - r_i^-$$

$$\min_{\mathbf{w}, r_i^+, r_i^-} \sum_i (r_i^+ + r_i^-) \quad s.t. \quad r_i^+ \geq 0, r_i^- \geq 0, \mathbf{w}^T \mathbf{x}_i + r_i^+ - r_i^- = t_i$$

- Note that with our definition above:

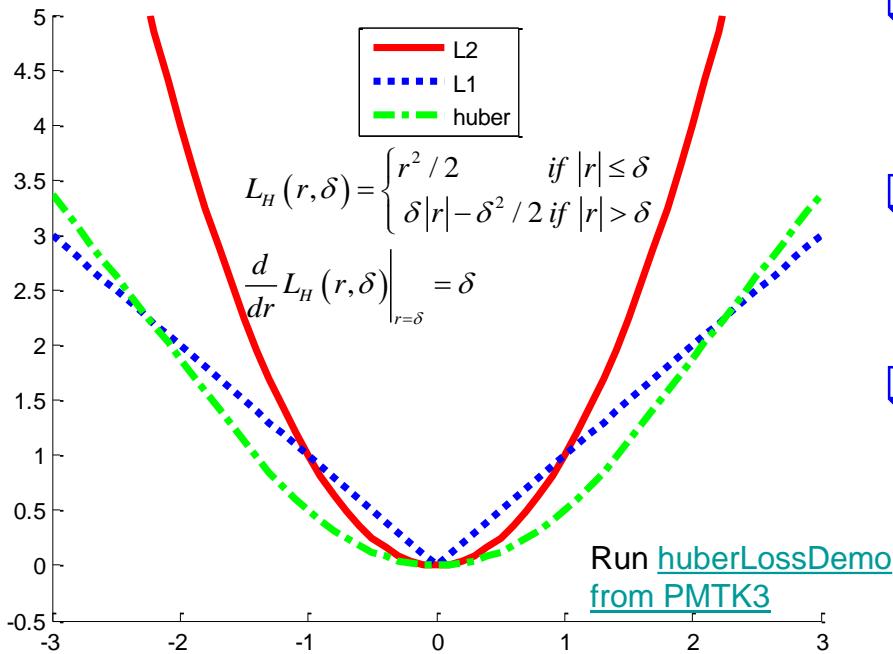
$$r_i^+ = \frac{1}{2}(r_i + |r_i|) = \begin{cases} r_i & \text{if } r_i \geq 0 \\ 0 & \text{if } r_i < 0 \end{cases}, \quad r_i^- = \frac{1}{2}(|r_i| - r_i) = \begin{cases} 0 & \text{if } r_i > 0 \\ -r_i & \text{if } r_i \leq 0 \end{cases}$$

$$|r_i| = r_i^+ + r_i^-$$

- Boyd, S. and L. Vandenberghe (2004). *Convex optimization*. Cambridge



# Huber Loss Function



- This is equivalent to  $L_2$  for errors that are smaller than  $\delta$ , and is equivalent to  $L_1$  for larger errors.
- This loss function is everywhere differentiable, using the fact that  $d/dr |r| = sign(r)$  if  $r \neq 0$ .
- The function is also  $C_1$  continuous, since the gradients of the two parts of the function match at  $r = \pm\delta$ .

- Optimizing the Huber loss is much faster than using the Laplace likelihood, since we can use standard optimization methods (quasi-Newton) rather than linear programming.
- The Huber method also has a probabilistic interpretation, although it is rather unnatural (Pontil et al. 1998).

- Pontil, M., S. Mukherjee, and F. Girosi (1998). [On the Noise Model of Support Vector Machine Regression](#). Technical report, MIT AI Lab.
- Huber, P. (1964). [Robust estimation of a location parameter](#). *Annals of Statistics* 53, 73-101.

# Regularized LS - Ridge Regression

- Consider the error function:

$$E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

data term + regularization term

- With the sum-of-squares error function and a quadratic regularizer, we get

$$\frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \lambda \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

- Specifically, setting the gradient with respect to  $\mathbf{w}$  to zero, and solving for  $\mathbf{w}$  as before, we obtain

$$\mathbf{w} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

- This is a trivial extension of the least-squares solution we encountered earlier (*Regularized Least Squares – Ridge Regression*)



# Regularized Least Squares

---

- Regularized solution:

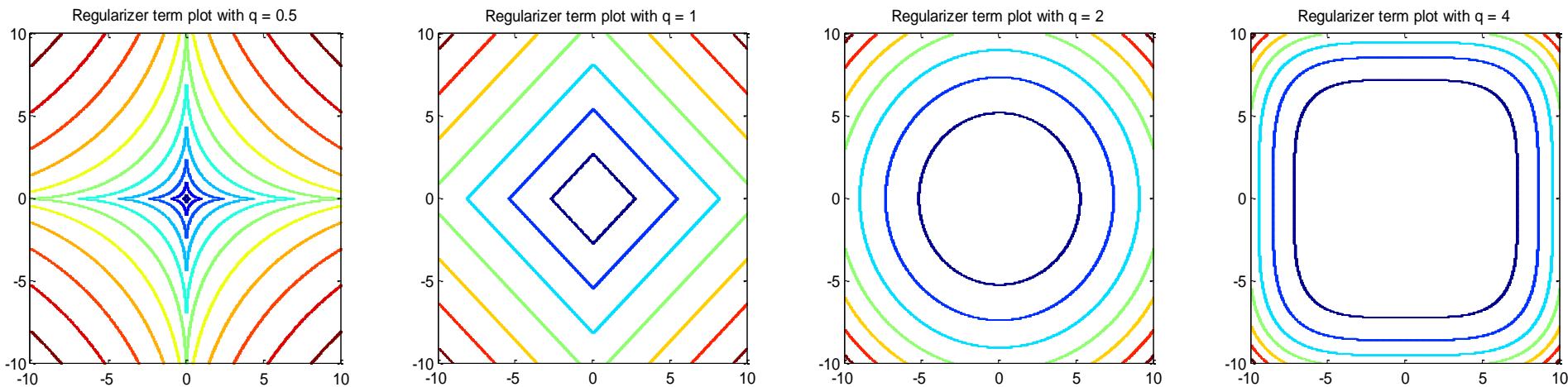
$$\mathbf{w} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

- Regularization limits the effective model complexity (the appropriate number of basis functions).
- This is replaced with the problem of finding a suitable value of the regularization coefficient  $\lambda$ .
- $\lambda$  controls how many non-zero  $w$ 's (i.e. basis functions) you have.

# Regularized Least Squares

- With a more general regularizer, we have

$$\frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \lambda \frac{1}{2} \sum_{j=1}^M |w_j|^q$$



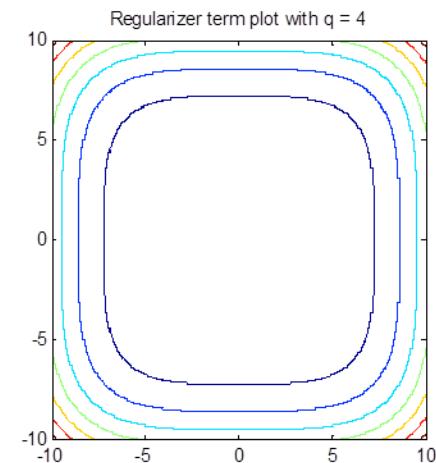
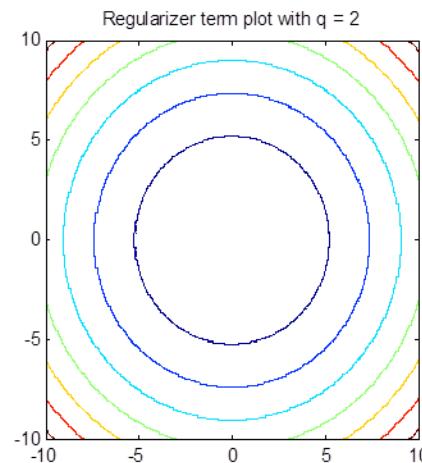
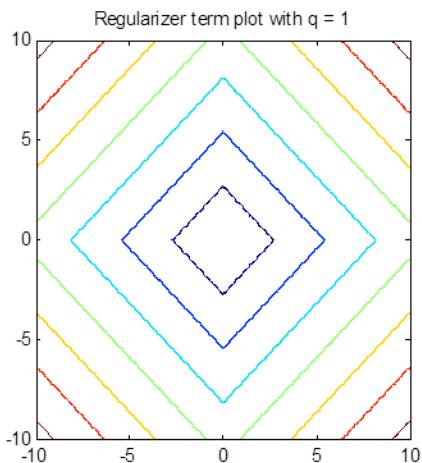
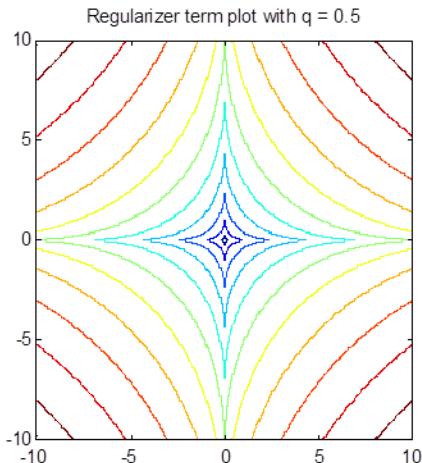
[MatLab code](#)

- $q = 1$  is known as the **Lasso regularizer**. These plots show only the regularizer term with  $\lambda = 0.7334$ .

# Regularized Least Squares

- With a more general regularizer, we have

$$\frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \lambda \frac{1}{2} \sum_{j=1}^M |w_j|^q$$

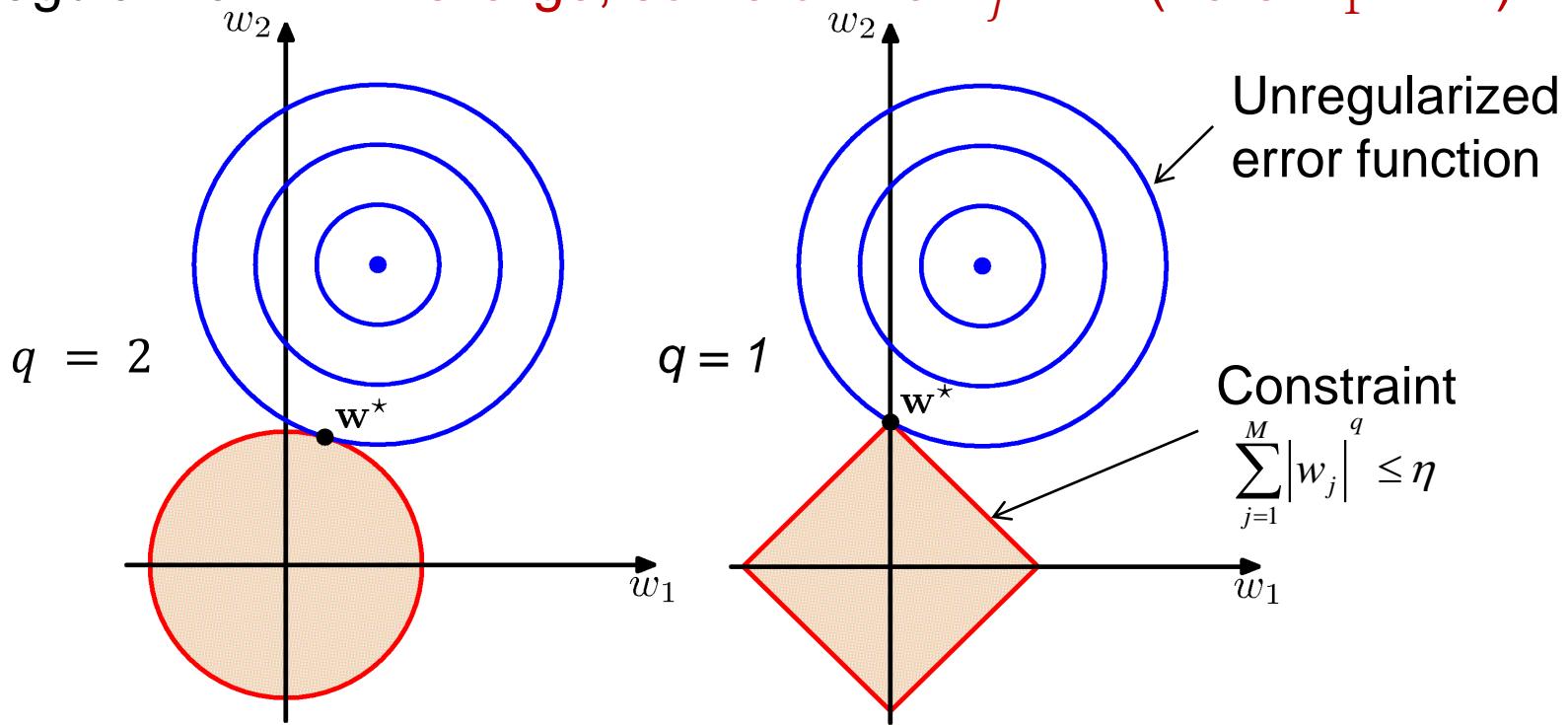


[MatLab code](#)

- $q = 2$  corresponds to the quadratic regularizer.

# Regularized Least Squares

- Lasso tends to generate sparser solutions than a quadratic regularizer – if  $\lambda$  is large, some of the  $w_j \rightarrow 0$  (here  $w_1 = 0$ ).



- Here, we consider that the regularized least squares solution is equivalent to minimizing the unregularized sum of squares with the constraint shown for some  $\eta$  (see proof next).

# Regularized Least Squares

- Let us write the constraint in the equivalent form:  $\frac{1}{2} \left( \sum_{j=1}^M |w_j|^q - \eta \right) \leq 0$
- This leads to the following Lagrangian function:

$$L(\mathbf{w}, \lambda) = \frac{1}{2} \sum_{n=1}^N \left( t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) \right)^2 + \frac{\lambda}{2} \left( \sum_{j=1}^M |w_j|^q - \eta \right)$$

- This is identical to our regularized least squares (RLS) in the dependence on  $\mathbf{w}$ .

$$\frac{1}{2} \sum_{n=1}^N \left( t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) \right)^2 + \lambda \frac{1}{2} \sum_{j=1}^M |w_j|^q \quad (*)$$

- For a particular  $\lambda > 0$ , let  $\mathbf{w}^*(\lambda)$  be the solution of the RLS in (\*).
- From the Kuhn-Tucker optimality conditions for  $L(\mathbf{w}, \lambda)$  we then see:

$$\eta = \sum_{j=1}^M |w_j^*|^q$$



# Kuhn-Tucker Optimality Conditions

- Consider the following constraint **minimization** problem:

$$\min_x f(\mathbf{x}), \text{subject to } g(\mathbf{x}) \geq 0$$

- This is equivalent as the minimization with respect to  $\mathbf{x}$  and  $\lambda$  of the following Lagrangian:

$$\min_{\mathbf{x}, \lambda} L(\mathbf{x}, \lambda) = f(\mathbf{x}) - \lambda g(\mathbf{x})$$

subject to the following (Kuhn-Tucker) conditions:

$$\lambda \geq 0, g(\mathbf{x}) \geq 0, \lambda g(\mathbf{x}) = 0$$

- Note for maximization problems, the Lagrangian should be modified as:  $L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$



# Multiple Outputs-Isotropic Covariance

- If we want to predict  $K > 1$  target variables, we use the same basis for all components of the target vector):

$$p(\mathbf{t} | \mathbf{x}, \mathbf{W}, \boldsymbol{\beta}) = \mathcal{N}(\mathbf{t} | \mathbf{y}(\mathbf{x}, \mathbf{W}), \boldsymbol{\beta}^{-1} \mathbf{I}) = \mathcal{N}(\mathbf{t} | \mathbf{W}^T \boldsymbol{\phi}(\mathbf{x}), \boldsymbol{\beta}^{-1} \mathbf{I})$$

$\mathbf{W}$  is an  $M \times K$  matrix and  $\mathbf{t}$  is  $K$  dimensional.

- Given observed inputs,  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , and targets,  $\mathbf{T} = \{\mathbf{t}_1, \dots, \mathbf{t}_N\}^T$  we obtain the log likelihood function

$$\ln p(\mathbf{T} | \mathbf{X}, \mathbf{W}, \boldsymbol{\beta}) = \sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{W}^T \boldsymbol{\phi}(\mathbf{x}_n), \boldsymbol{\beta}^{-1} \mathbf{I}) = \frac{NK}{2} \ln \frac{\boldsymbol{\beta}}{2\pi} - \frac{\boldsymbol{\beta}}{2} \sum_{n=1}^N \| \mathbf{t}_n - \mathbf{W}^T \boldsymbol{\phi}(\mathbf{x}_n) \|^2$$

# **K-Independent Regression Problems**

$$\ln p(\mathbf{T} | \mathbf{X}, \mathbf{W}, \beta) = \frac{NK}{2} \ln \frac{\beta}{2\pi} - \frac{\beta}{2} \sum_{n=1}^N \left\| \mathbf{t}_n - \mathbf{W}^T \boldsymbol{\phi}(\mathbf{x}_n) \right\|^2$$

- As before, we can maximize this function with respect to  $\mathbf{W}$ , giving

$$\mathbf{W}_{ML} = \underbrace{\left( \boldsymbol{\Phi}^T \boldsymbol{\Phi} \right)^{-1}}_{M \times M} \underbrace{\boldsymbol{\Phi}^T}_{M \times N} \underbrace{\mathbf{T}}_{N \times K}$$

- If we examine this result for each target variable  $t_k$ , we have (take the  $k^{\text{th}}$  column of  $\mathbf{W}$  and  $\mathbf{T}$ ):

$$\mathbf{w}_{k_{ML}} = \left( \boldsymbol{\Phi}^T \boldsymbol{\Phi} \right)^{-1} \boldsymbol{\Phi}^T \mathbf{t}_k = \boldsymbol{\Phi}^\dagger \mathbf{t}_k$$

which is identical with the single output case (so **there is decoupling between the target variables**)

- As expected, we obtain  $K$  –**independent regression problems.**



# Multiple Outputs - Full Covariance

- Let us repeat the earlier formulation but with covariance matrix  $\Sigma$ . If we want to predict  $K > 1$  target variables, we use the same basis for all components of the target vector):

$$p(\mathbf{t} | \mathbf{x}, \mathbf{W}, \Sigma) = \mathcal{N}(\mathbf{t} | \mathbf{y}(\mathbf{x}, \mathbf{W}), \Sigma) = \mathcal{N}(\mathbf{t} | \mathbf{W}^T \phi(\mathbf{x}), \Sigma)$$

where  $\mathbf{W}$  is an  $M \times K$  matrix of parameters

- Given observed inputs,  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , and targets,  $\mathbf{T} = \{\mathbf{t}_1, \dots, \mathbf{t}_N\}^T$  we obtain the log likelihood function  $\ln p(\mathbf{T} | \mathbf{X}, \mathbf{W}, \Sigma)$  :

$$\sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{W}^T \phi(x_n), \Sigma) = -\frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}^T \phi(\mathbf{x}_n))^T \Sigma^{-1} (\mathbf{t}_n - \mathbf{W}^T \phi(\mathbf{x}_n))$$

# Multiple Outputs - Full Covariance

$$\ln p(\mathbf{T} | \mathbf{X}, \mathbf{W}, \Sigma) = -\frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}^T \boldsymbol{\phi}(\mathbf{x}_n))^T \Sigma^{-1} (\mathbf{t}_n - \mathbf{W}^T \boldsymbol{\phi}(\mathbf{x}_n))$$

- As before, we maximize this function with respect to  $\mathbf{W}$ ,

$$0 = -\sum_{n=1}^N \Sigma^{-1} (\mathbf{t}_n - \mathbf{W}^T \boldsymbol{\phi}(\mathbf{x}_n)) \boldsymbol{\phi}(\mathbf{x}_n)^T \Rightarrow \mathbf{W}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{T}$$

- For the ML estimate for  $\Sigma$ , use the result for the MLE of the covariance of a multivariate Gaussian:

$$\Sigma = \frac{1}{N} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}_{ML}^T \boldsymbol{\phi}(\mathbf{x}_n)) (\mathbf{t}_n - \mathbf{W}_{ML}^T \boldsymbol{\phi}(\mathbf{x}_n))^T$$

- Note that each column of  $\mathbf{W}_{ML}$  is of the form

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

seen for isotropic noise distribution and is independent of  $\Sigma$ !

# The Bias-Variance Decomposition

---

- MLE (i.e. least squares) leads to severe over-fitting if complex models are trained using data sets of limited size.
- Over-fitting occurs whenever the number of basis functions is large (i.e. for complex models) and the training data set is of limited size.
- Limiting the number of basis functions limits the flexibility of the model.
- Regularization controls over-fitting but one needs to determine  $\lambda$ .
- Over-fitting is property of MLE and does not arise when we marginalize over parameters in a Bayesian setting.
- Before returning to a Bayesian setting, we will discuss the bias-variance tradeoff (a frequentist viewpoint of model complexity).



# **Loss Function and the Regression Function**

---

- Recall the regression loss-function

$$L(t, y(\mathbf{x})) = (y(\mathbf{x}) - t)^2$$

- The decision problem is to minimize the expected loss:

$$\mathbb{E}[L] = \int \int (y(\mathbf{x}) - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

- The solution of this minimization problem known as the regression function is:

$$y(\mathbf{x}) = \int t p(t|\mathbf{x}) dt = \mathbb{E}_t[t|\mathbf{x}]$$

i.e. the average of  $t$  conditioned on  $\mathbf{x}$ .

- A useful expression was derived in an earlier lecture:

$$\mathbb{E}[L] = \int (y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}])^2 p(\mathbf{x}) d\mathbf{x} + \int \int (\mathbb{E}[t|\mathbf{x}] - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

# The Bias-Variance Decomposition

- We can write the expected squared loss as:

$$\mathbb{E}[L] = \int (y(\mathbf{x}) - h(\mathbf{x}))^2 p(\mathbf{x}) d\mathbf{x} + \iint (h(\mathbf{x}) - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

where  $h(\mathbf{x})$  is the conditional expectation which is given by

$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int t p(t|\mathbf{x}) dt$$

- Recall that the second term, which is independent of  $y(\mathbf{x})$ , arises from the intrinsic noise on the data and represents the minimum achievable value of the expected loss.
- If we model  $h(\mathbf{x})$  using a parametric function  $y(\mathbf{x}, \mathbf{w})$  governed by  $\mathbf{w}$ , then from a Bayesian perspective the uncertainty in our model is expressed through a posterior distribution over  $\mathbf{w}$ .



# The Bias-Variance Decomposition

$$\mathbb{E}[L] = \int (y(x) - h(x))^2 p(x) dx + \iint (h(x) - t)^2 p(x, t) dx dt$$

- A frequentist treatment involves making a point estimate of  $w$  based on the data set  $\mathcal{D}$ , and interpret the uncertainty of this estimate as follows.
- Suppose we had a large number of data sets each of size  $N$  and each drawn independently from  $p(x, t)$ .
- For any given data set  $\mathcal{D}$ , we run our learning algorithm and obtain a prediction function  $y(x; \mathcal{D})$ . Different data sets from the ensemble give different functions and values of the squared loss.
- The performance of the learning algorithm is then assessed by taking the average over this ensemble of data sets.

# The Bias-Variance Decomposition

- For any given data set  $\mathcal{D}$ , we can run our learning algorithm and obtain a prediction function  $y(\mathbf{x}; \mathcal{D})$ .

$$\begin{aligned}\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2 &= \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] + \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ &= \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2 + \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 + \\ &\quad 2\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}\end{aligned}$$

- Here  $\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]$  is the average of our prediction function over all data sets.
- Take the expectation of this expression with respect to  $\mathcal{D}$  and note that the final term will vanish, giving

$$\mathbb{E}_{\mathcal{D}}[\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2] = \underbrace{\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2}_{(Bias)^2} + \underbrace{\mathbb{E}_{\mathcal{D}}[\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}]^2}_{Variance}$$

- Recall that  $h(\mathbf{x})$  is the desired regression function.

# The Bias-Variance Decomposition

$$\mathbb{E}[L] = \int (y(\mathbf{x}) - h(\mathbf{x}))^2 p(\mathbf{x}) d\mathbf{x} + \iint (h(\mathbf{x}) - t)^2 p(\mathbf{x}, t) dx dt$$

- So far, we have considered a single input value  $\mathbf{x}$ . If we substitute this expansion back into the expected squared loss function shown above, we obtain the following decomposition of the expected squared loss

Expected loss = (bias)<sup>2</sup> + variance + noise

where

$$bias^2 = \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x}$$

$$variance = \int \mathbb{E}_{\mathcal{D}}[\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) d\mathbf{x}$$

$$noise = \iint (h(\mathbf{x}) - t)^2 p(\mathbf{x}, t) dx dt$$



# The Bias-Variance Decomposition

Expected loss = (bias)<sup>2</sup> + variance + noise

$$\text{bias}^2 = \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x}$$

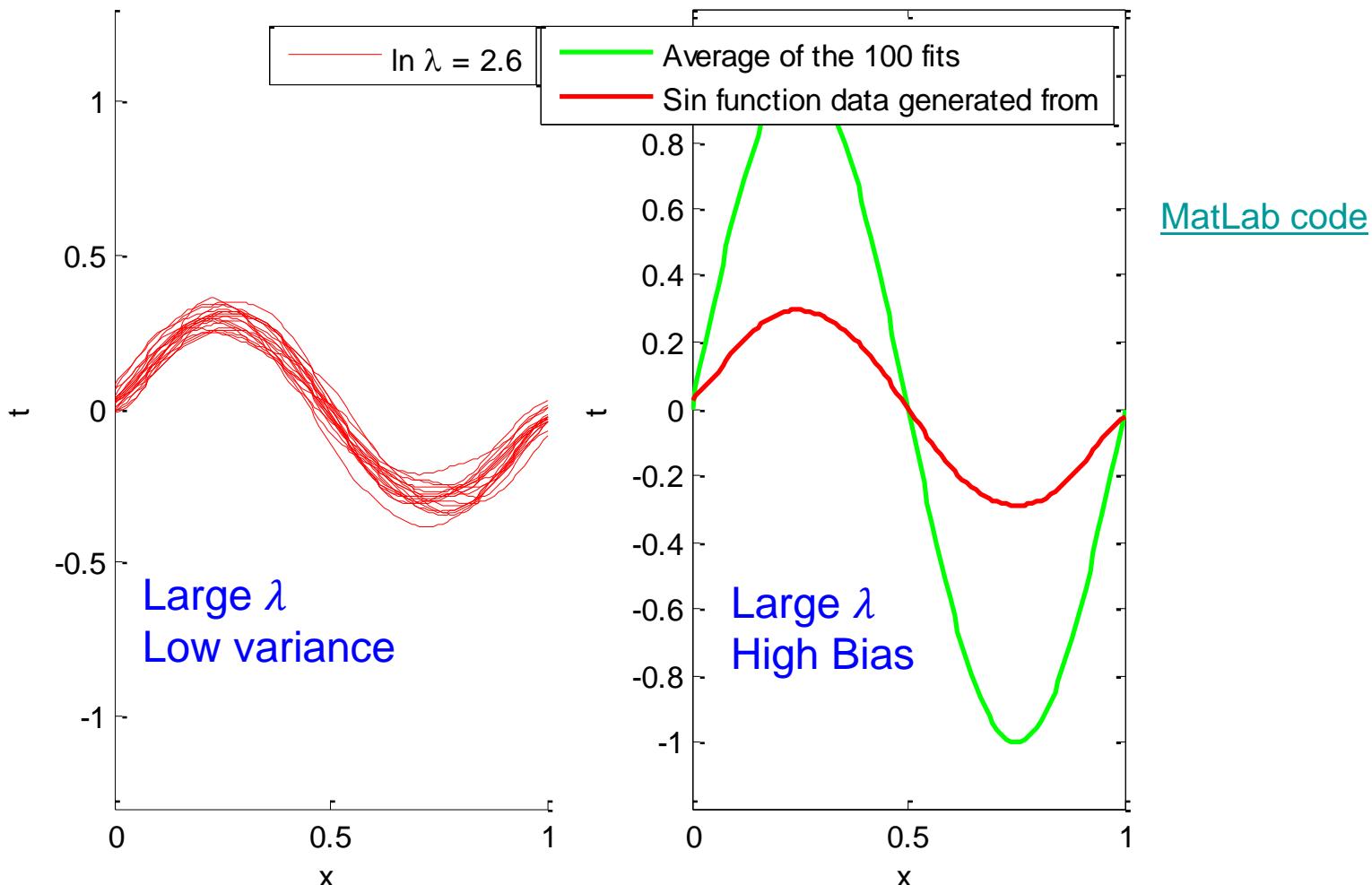
$$\text{variance} = \int \mathbb{E}_{\mathcal{D}}[\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) d\mathbf{x}$$

$$\text{noise} = \iint (h(\mathbf{x}) - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

- There is a tradeoff between bias and variance
  - flexible models have low bias and high variance
  - rigid models have high bias and low variance
- The bias-variance decomposition provides insights in model complexity but is of limited use since several data sets  $\mathcal{D}$  are needed.

# The Bias-Variance Decomposition

$L = 100$  data sets, each with  $N = 25$  data points. 24 Gaussian basis functions, # of parameters  $M = 25$

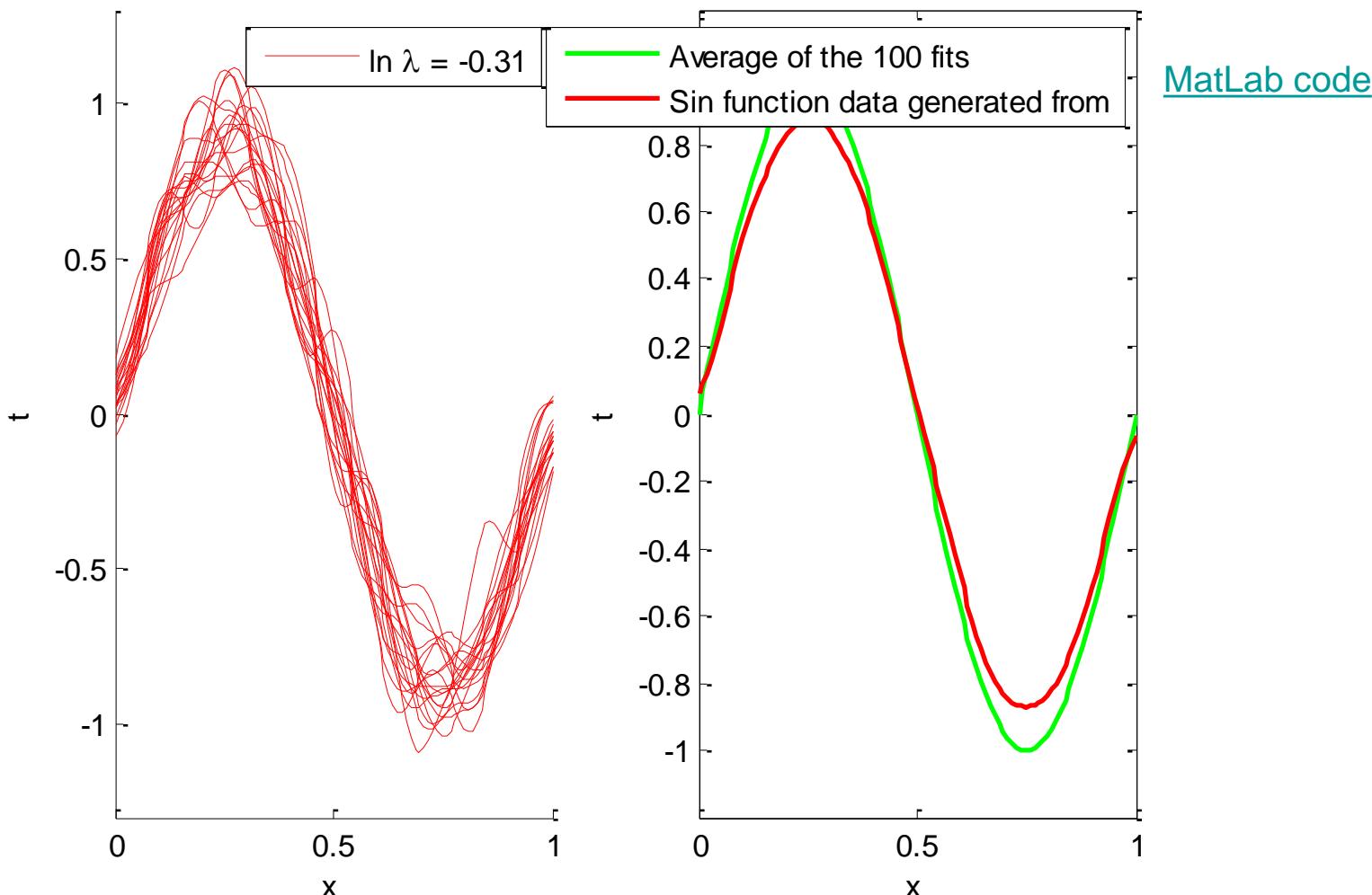


- Dependence of bias and variance on model complexity, governed by a regularization parameter  $\lambda$ , using the sinusoidal data set. On the left only 20 of the 100 fits are shown.



# The Bias-Variance Decomposition

$L = 100$  data sets, each with  $N = 25$  data points. 24 Gaussian basis functions, # of parameters  $M = 25$



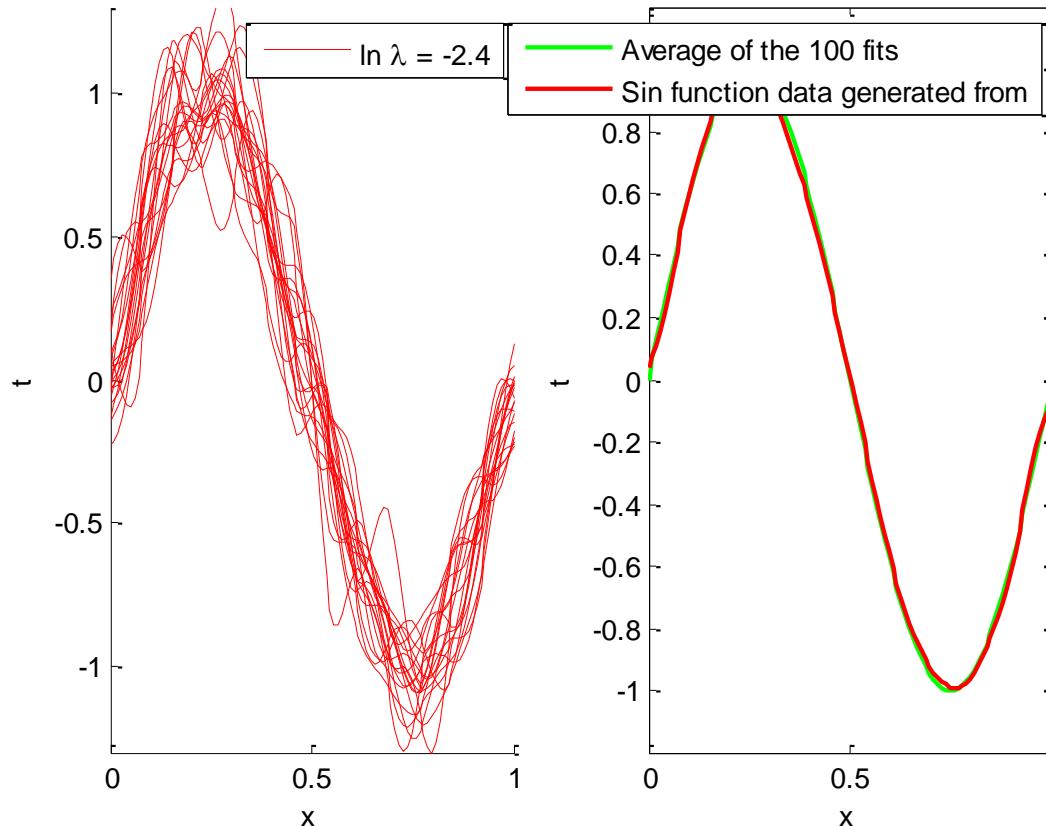
[MatLab code](#)

- Dependence of bias and variance on model complexity, governed by a regularization parameter  $\lambda$ , using the sinusoidal data set. On the left only 20 of the 100 fits are shown.

# The Bias-Variance Decomposition

$L = 100$  data sets, each with  $N = 25$  data points. 24 Gaussian basis functions, # of parameters  $M = 25$

Lowest  $\lambda$   
Large  
variance  
Low Bias



[MatLab code](#)

- Note averaging of many solutions of the complex model ( $M=25$ ) obtained from different data sets gives very good solution.
- Weighted averaging is also done in Bayesian approach but there it is with respect to the posterior distribution of parameters!

# Trade-off quantities

---

- The average prediction over  $L$  data sets is estimated from

$$\bar{y}(x) = \frac{1}{L} \sum_{\ell=1}^L y^{(\ell)}(x)$$

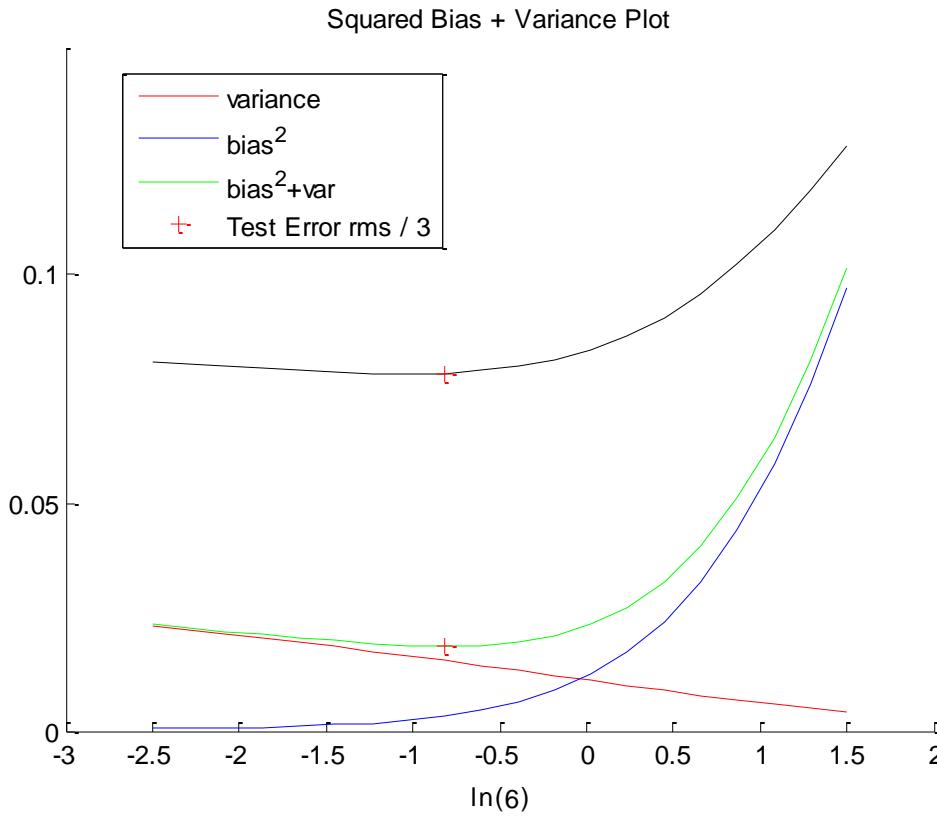
- The integrated squared bias and integrated variance are then given by

$$bias^2 = \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} \Rightarrow (bias)^2 = \frac{1}{N} \sum_{n=1}^N \{\bar{y}(\mathbf{x}_n) - h(\mathbf{x}_n)\}^2$$

$$variance = \int \mathbb{E}_{\mathcal{D}}[\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) d\mathbf{x} \Rightarrow$$

$$variance = \frac{1}{N} \sum_{n=1}^N \frac{1}{L} \sum_{\ell=1}^L \{y^{(\ell)}(\mathbf{x}_n) - \bar{y}(\mathbf{x}_n)\}^2$$

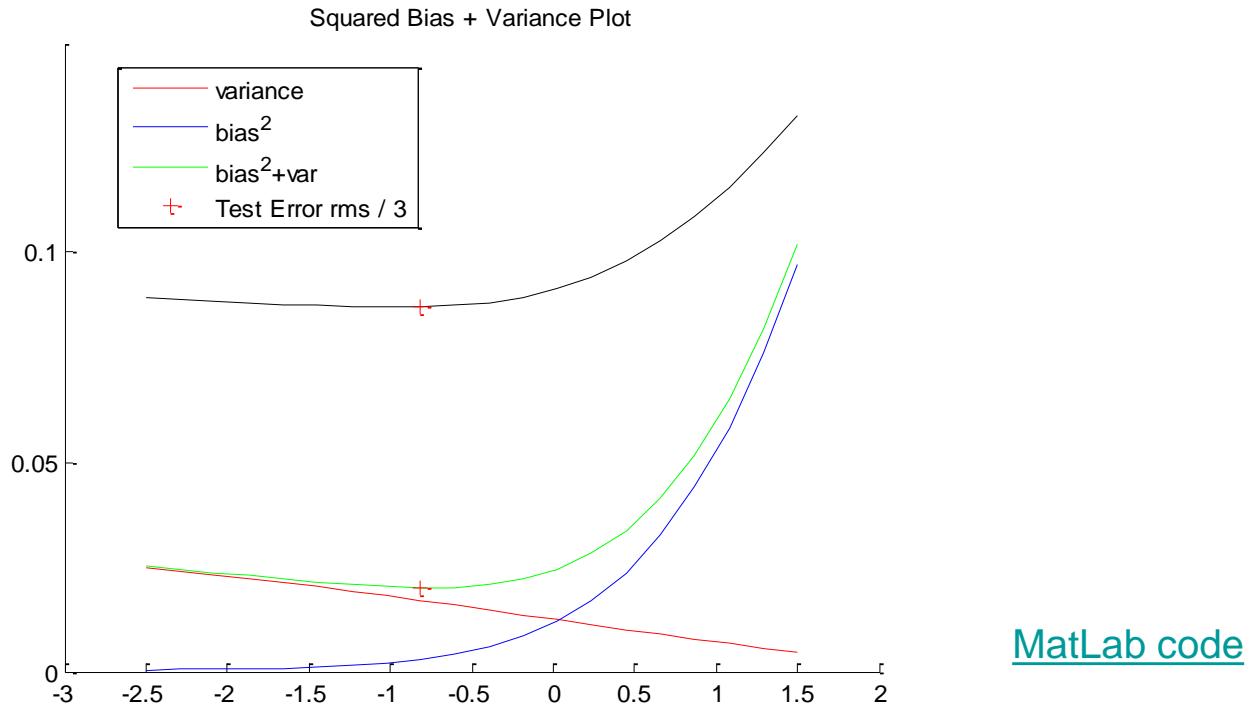
# The Bias-Variance Decomposition



[MatLab code](#)

- ❑ Note that small regularization parameter  $\lambda$  allows the model to be finely tuned to the noise in each individual data set leading to small bias but large variance.
- ❑ Conversely, large  $\lambda$  forces all  $w$ 's to go towards zero leading to large bias but small variance.

# The Bias-Variance Decomposition



- Plot of squared bias and variance together with their sum. Also shown is the test set rms squared error for a test data set size of 1000 points.
- The minimum value of  $(\text{bias})^2 + \text{variance}$  occurs around  $\ln \lambda = -0.31$ , which is near to the value that gives the minimum error on the test data.

# MAP and Regularized Least Squares

- Now let us take a step towards a Bayesian approach and introduce a prior distribution over the coefficients  $\mathbf{w}$ .

Consider a Gaussian distribution

$$p(\mathbf{w} | \alpha) = \mathcal{N}(\mathbf{w} | 0, \alpha^{-1} I) = \left( \frac{\alpha}{2\pi} \right)^{M/2} \exp \left\{ -\frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \right\}$$

where  $\alpha$  is the precision of the distribution, and  $M$  is the total number of elements in the vector  $\mathbf{w}$ .\*

- Using Bayes' theorem:

$$p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) p(\mathbf{w} | \alpha)$$

\* Note that one should not penalize the bias term  $\mathbf{w}_0$  as it does not contribute to overfitting. We will discuss on a follow up lecture how to address this by treating the bias term separately by working with “centered data” (both responses  $\mathbf{t}_c = \mathbf{t} - \bar{\mathbf{t}} = \mathbf{t} - \mathbf{1}_N \bar{\mathbf{t}}$ ,  $\bar{\mathbf{t}} \equiv \sum_{i=1}^N t_i / N$ , and input variables s.t.  $y(x_n, \mathbf{w}) = \boldsymbol{\phi}_c(x_n)^T \mathbf{w}$ ,  $\boldsymbol{\Phi}_c = (\boldsymbol{\varphi}_{1c} \quad \boldsymbol{\varphi}_{2c} \quad \dots \quad \boldsymbol{\varphi}_{Mc})$ ,  $\boldsymbol{\varphi}_{ic} = \boldsymbol{\varphi}_i - \frac{1}{N} \sum_{n=1}^N \boldsymbol{\varphi}_i(x_n) = 0$ ,  $i = 1, \dots, M-1$ ). Here  $\boldsymbol{\phi}^T$  are rows and  $\boldsymbol{\varphi}$  columns of the  $N \times M$  design matrix  $\boldsymbol{\Phi}$ .



# MAP Estimate

---

- We can determine  $w$  (point estimate) by finding the most probable value of  $w$  given the data, i.e. maximizing the posterior.
- This technique is called *maximum posterior (MAP)*.
- The maximum of the posterior is given by the minimum of

$$\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2 + \frac{\alpha}{2} w^T w$$

- Note that the MAP point estimate is equivalent to regularized sum of squares error function with regularization parameter

$$\lambda = \frac{\alpha}{\beta} = \frac{\text{precision of prior}}{\text{precision in the data}}$$



# Posterior Distribution

$$p(\mathbf{w} | \alpha) = \left( \frac{\alpha}{2\pi} \right)^{M/2} \exp \left\{ -\frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \right\}, p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) \propto \exp \left\{ -\frac{\beta}{2} \sum_{n=1}^N (\boldsymbol{\phi}(x_n)^T \mathbf{w} - t_n)^2 \right\} \Rightarrow$$

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) \propto \exp \left\{ -\frac{1}{2} \mathbf{w}^T \sum_{n=1}^N \beta \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T \mathbf{w} + \beta \mathbf{w}^T \sum_{n=1}^N t_n \boldsymbol{\phi}(x_n) \right\}$$

Quadratic in  $\mathbf{w}$

- We now have the product of two Gaussians and the posterior is easily computed as:

$$p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \alpha, \beta) \propto$$

$$\exp \left( -\frac{1}{2} \mathbf{w}^T \alpha \mathbf{I}_{M \times M} \mathbf{w} - \frac{1}{2} \mathbf{w}^T \sum_{n=1}^N \beta \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T \mathbf{w} + \beta \mathbf{w}^T \sum_{n=1}^N t_n \boldsymbol{\phi}(x_n) \right)$$

$$\propto \mathcal{N} \left( \mathbf{w} | \beta \mathbf{S}_N \sum_{n=1}^N t_n \boldsymbol{\phi}(x_n), \mathbf{S}_N \right), \mathbf{S}_N^{-1} = \alpha \mathbf{I} + \sum_{n=1}^N \beta \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T$$

Completing the square  
 $\downarrow$   

$$-\frac{1}{2} (\mathbf{w} - \boldsymbol{\mu})^T \mathbf{S}_N^{-1} (\mathbf{w} - \boldsymbol{\mu})$$

E.g. for polynomial regression:  $\boldsymbol{\phi}(x_n) = \begin{bmatrix} 1 \\ x_n \\ x_n^2 \\ \vdots \\ x_n^{M-1} \end{bmatrix}, \boldsymbol{\phi}(x)^T = \{1 \quad x \quad x^2 \quad \dots \quad x^{M-1}\}, \mathbf{I} = \text{unit matrix } M \times M$

# Posterior Distribution

---

- For a prior

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I})$$

the corresponding posterior distribution over  $\mathbf{w}$  is then given by

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N)$$

$$\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t}$$

$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi$$



# General Gaussian as a Prior $p(\mathbf{w})$

- Assume additive Gaussian noise with known precision  $\beta$ . The likelihood function  $p(\mathbf{t}|\mathbf{w})$  is the exponential of a quadratic function of  $\mathbf{w}$

$$p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}) = \left( \frac{\beta}{2\pi} \right)^{N/2} \exp \left( -\frac{\beta}{2} \sum_{n=1}^N \{ t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) \}^2 \right)$$

and its **conjugate prior** is Gaussian:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$$

- Combining this with the likelihood and using results for marginal and conditional Gaussian distributions, gives **the posterior**

$$p(\mathbf{w} | \mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$$

where

$$\mathbf{m}_N = \mathbf{S}_N \left( \mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \boldsymbol{\Phi}^T \mathbf{t} \right)$$

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

# Posterior Distribution: Derivation

$$p(\mathbf{w} | \mathbf{m}_0, S_0) \propto \exp \left\{ -\frac{1}{2} (\mathbf{w} - \mathbf{m}_0)^T S_0^{-1} (\mathbf{w} - \mathbf{m}_0) \right\},$$

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) \propto \exp \left\{ -\frac{\beta}{2} \sum_{n=1}^N (\phi(x_n)^T \mathbf{w} - t_n)^2 \right\} \Rightarrow$$

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) \propto \exp \left\{ -\frac{1}{2} \mathbf{w}^T \sum_{n=1}^N \beta \phi(x_n) \phi(x_n)^T \mathbf{w} + \beta \sum_{n=1}^N t_n \phi(x_n) \mathbf{w} \right\}$$

- We now have the product of two Gaussians and the posterior is easily computed as:

$$p(\mathbf{w} | \mathbf{x}, \mathbf{t}, S_0, \beta) \propto$$

$$\exp \left( -\frac{1}{2} \mathbf{w}^T S_0^{-1} \mathbf{w} - \mathbf{w}^T S_0^{-1} \mathbf{m}_0 - \frac{1}{2} \mathbf{w}^T \sum_{n=1}^N \beta \phi(x_n) \phi(x_n)^T \mathbf{w} + \beta \mathbf{w}^T \sum_{n=1}^N t_n \phi(x_n) \right)$$

$$\propto \mathcal{N} \left( \mathbf{w} / \underbrace{S_N \left( S_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t} \right)}_{\mathbf{m}_N}, S_N \right), S_N^{-1} = S_0^{-1} + \sum_{n=1}^N \beta \phi(x_n) \phi(x_n)^T = S_0^{-1} + \beta \Phi^T \Phi$$

Complete the square in  $\mathbf{w}$



# Sequential Posterior Calculation

- Note that because the posterior distribution is Gaussian, its posterior mode coincides with its mean.

$$p(\mathbf{w} | \mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N) \quad \mathbf{m}_N = \mathbf{S}_N^{-1} (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t}) \\ \mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \Phi^T \Phi$$

$$\mathbf{w}_{MAP} = \mathbf{m}_N$$

- The above expressions for the posterior mean and variance can also be written for a **sequential calculation** (we already have observed  $N$  data points and now considering an additional data point  $(\mathbf{x}_{N+1}, t_{N+1})$ ). In this case, we have:

$$p(\mathbf{w} | \mathbf{t}_{N+1}, \mathbf{x}_{N+1}, \mathbf{m}_N, \mathbf{S}_N) = \mathcal{N}(\mathbf{w} | \mathbf{m}_{N+1}, \mathbf{S}_{N+1})$$

$$\mathbf{m}_{N+1} = \mathbf{S}_{N+1}^{-1} (\mathbf{S}_N^{-1} \mathbf{m}_N + \beta \phi_{n+1} t_{n+1}) \\ \mathbf{S}_{N+1}^{-1} = \mathbf{S}_N^{-1} + \beta \phi_{n+1} \phi_{n+1}^T$$

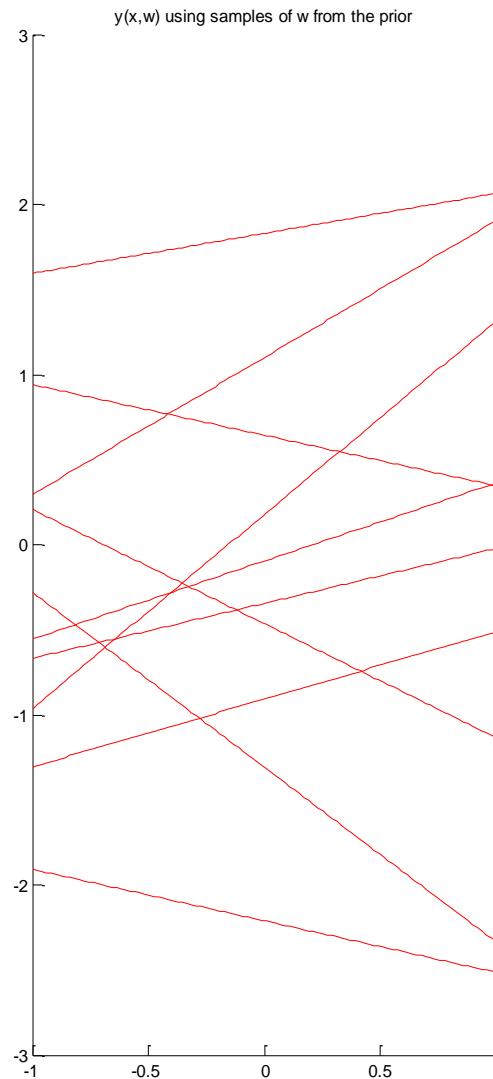
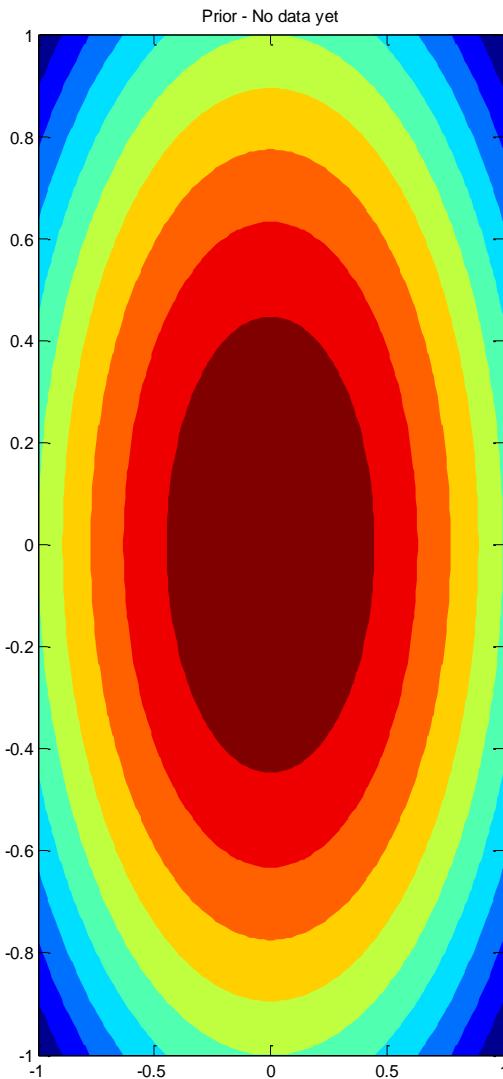
# Bayesian Regression: Example

---

- We generate synthetic data from the function  $f(x, a) = a_0 + a_1 x$  with parameter values  $a_0 = -0.3$  and  $a_1 = 0.5$  by first choosing values of  $x_n$  from the uniform distribution  $\mathcal{U}(x| -1, 1)$ , then evaluating  $f(x_n, a)$ , and finally adding Gaussian noise with standard deviation of 0.2 to obtain the target values  $t_n$ .
- We assume  $\beta = (1/0.2)^2 = 25$  and  $\alpha = 2.0$ .
- We perform Bayesian inference sequentially – one point at a time – so the posterior at each level becomes the new prior.
- We show results after 1, 2 and 22 points have been collected.
- The results include the likelihood contours (for 1 point), the posterior and samples of the regression function from the posterior.



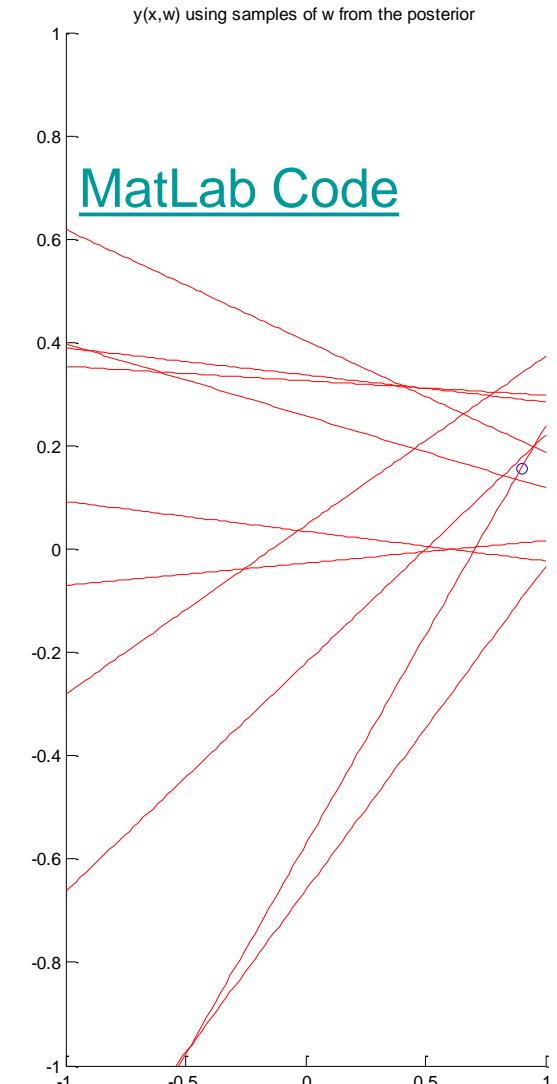
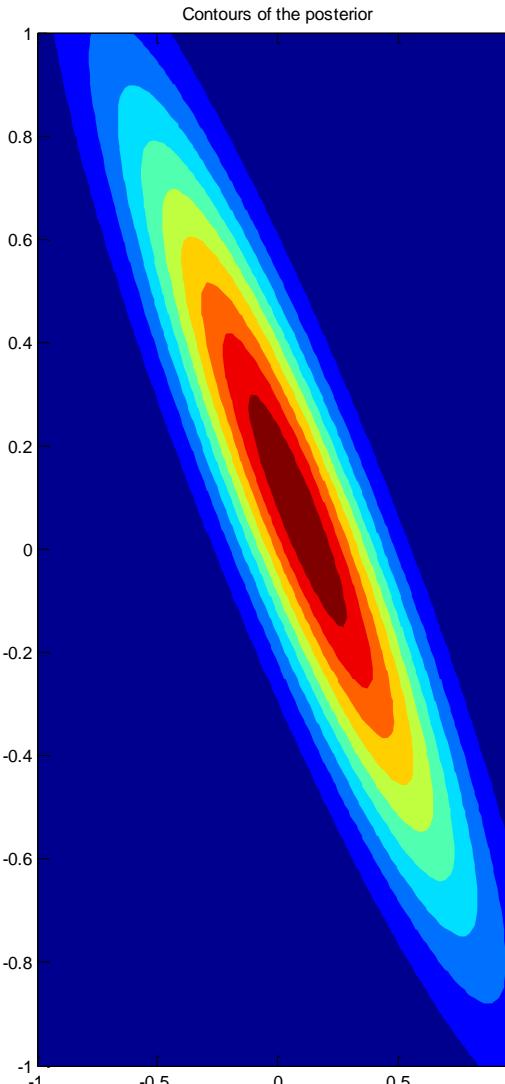
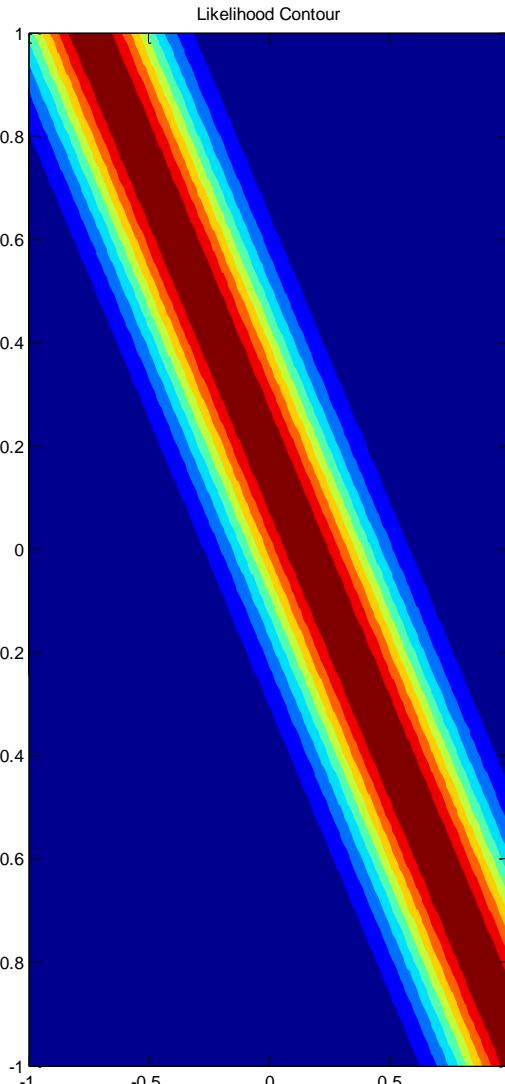
# Bayesian Regression: Example



[MatLab Code](#)



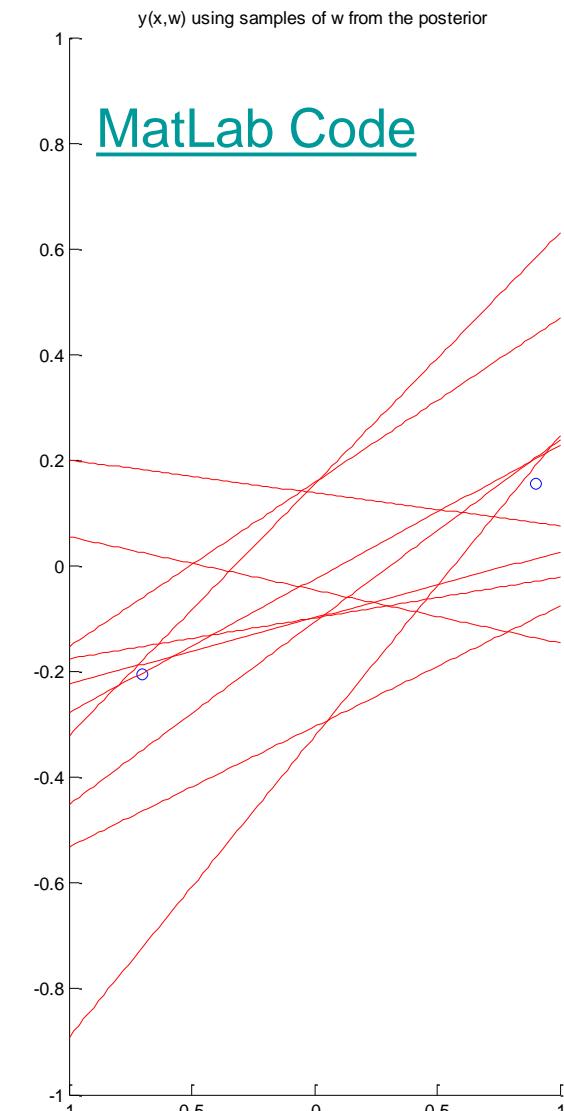
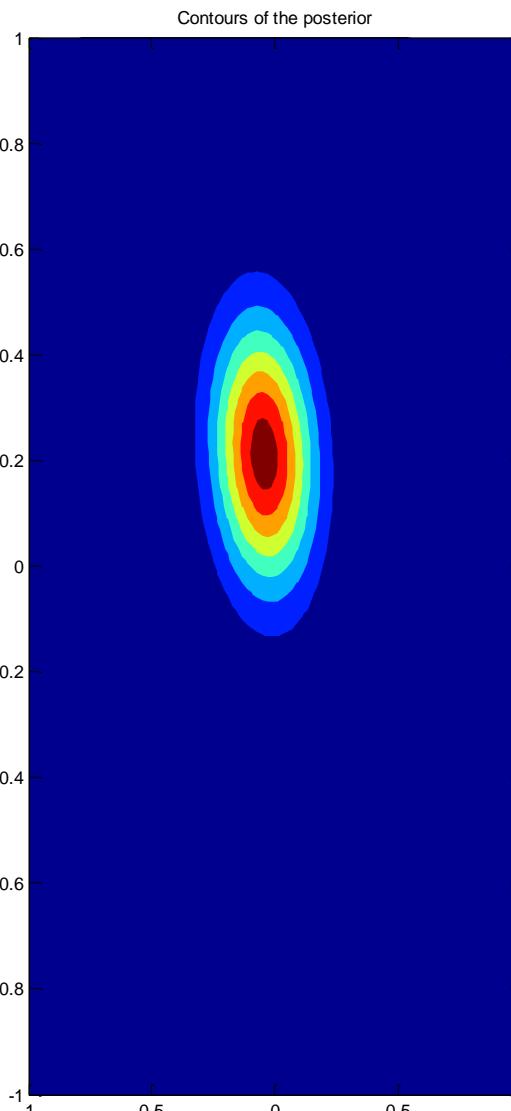
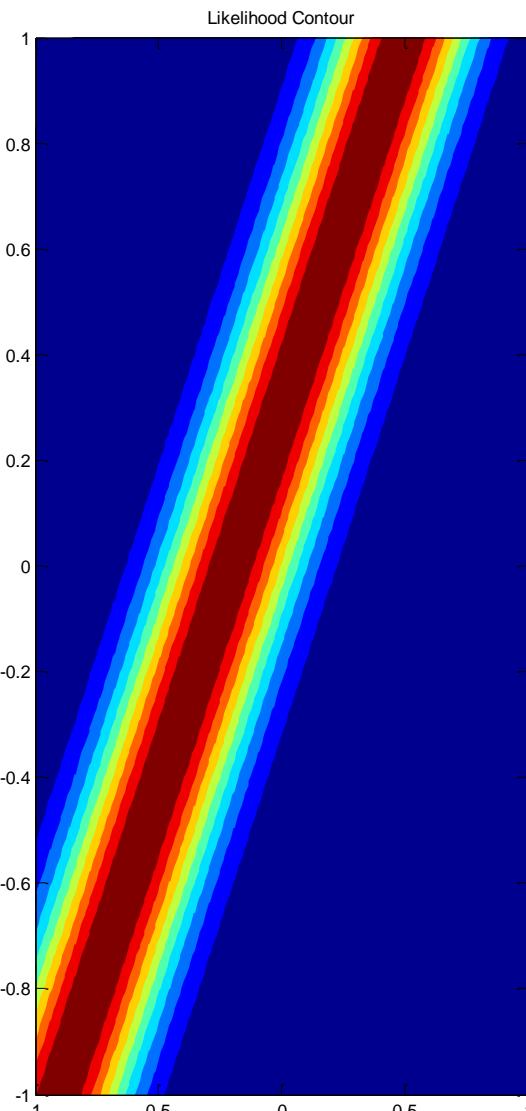
# Example: One Data Point Collected



Note that the regression lines pass close to the data point (shown with a circle)



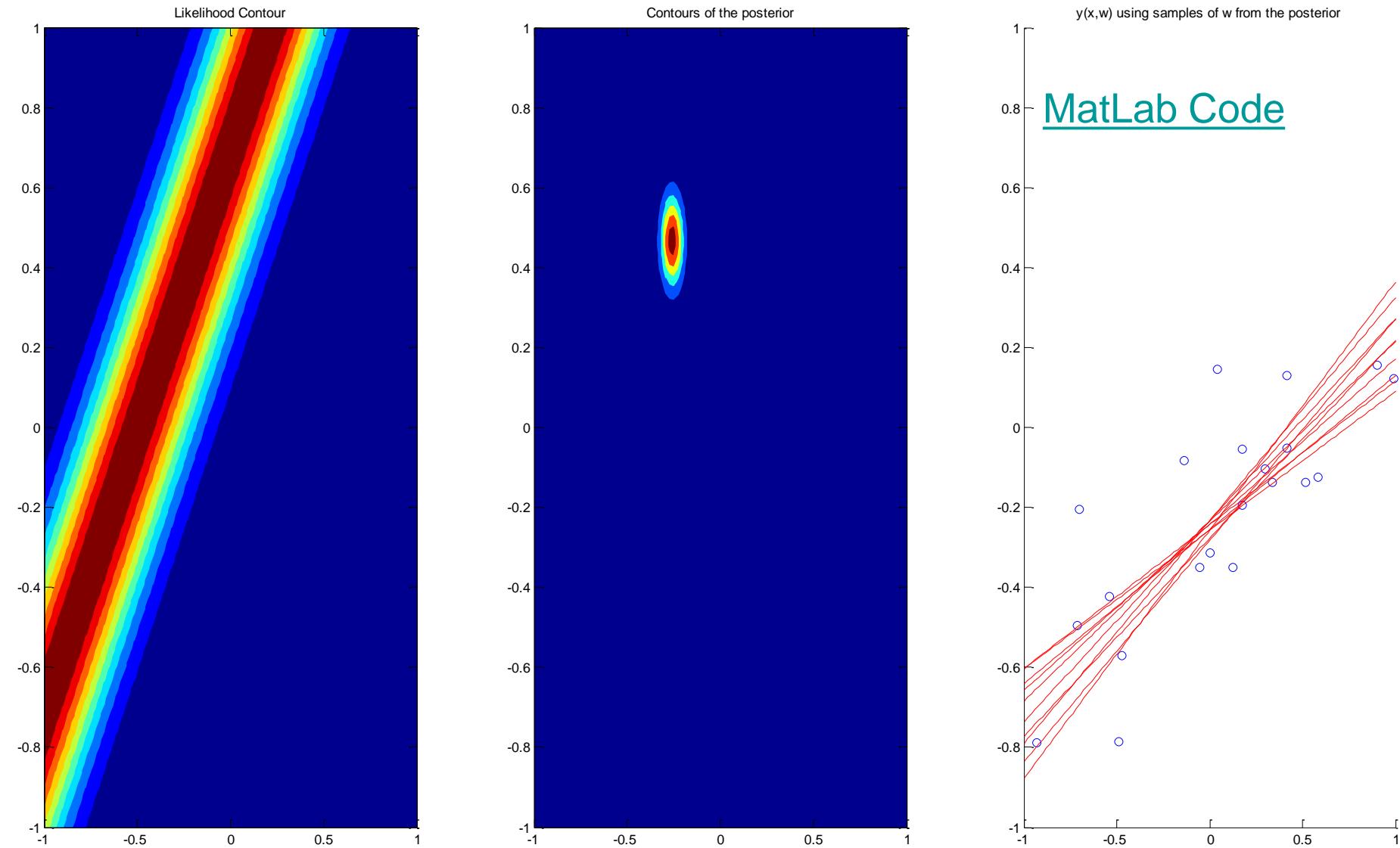
# Example: 2<sup>nd</sup> Data Point Collected



Note that the regression lines now pass close to both data points



# Example: 22 Data Points Collected

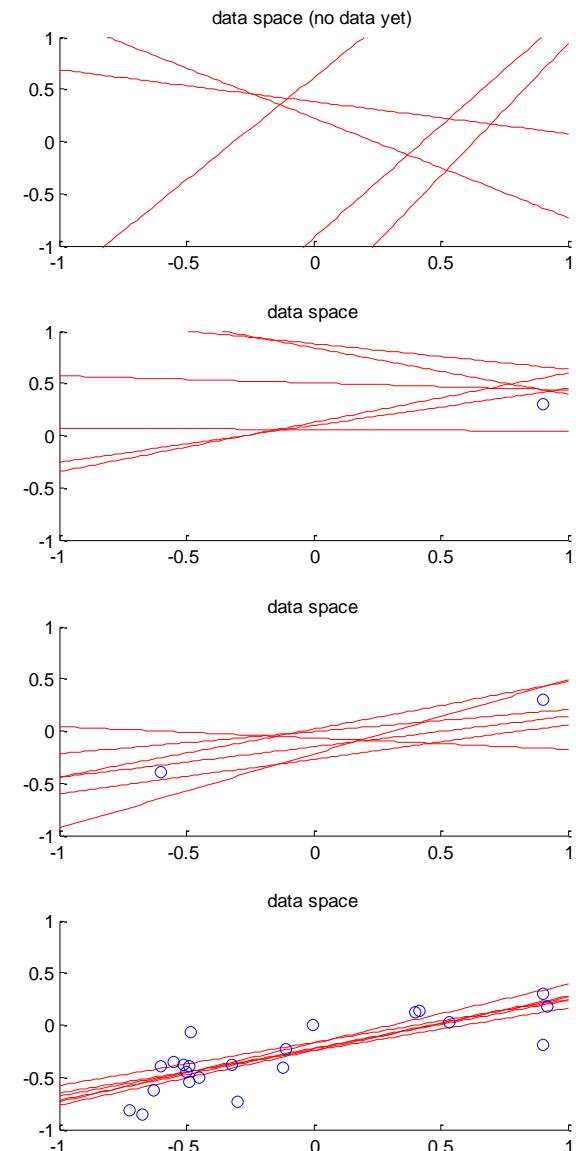
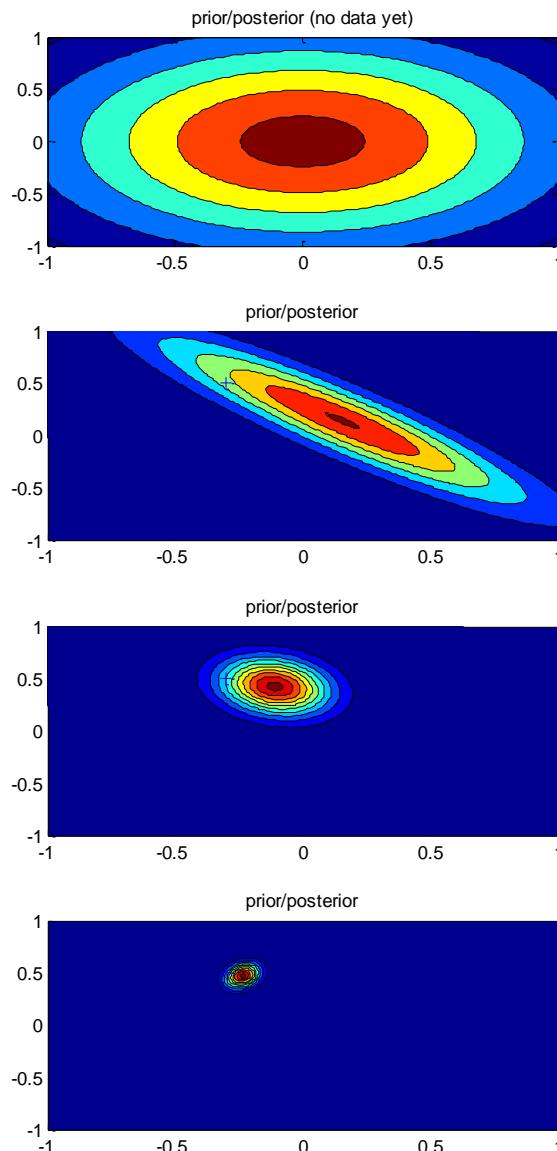
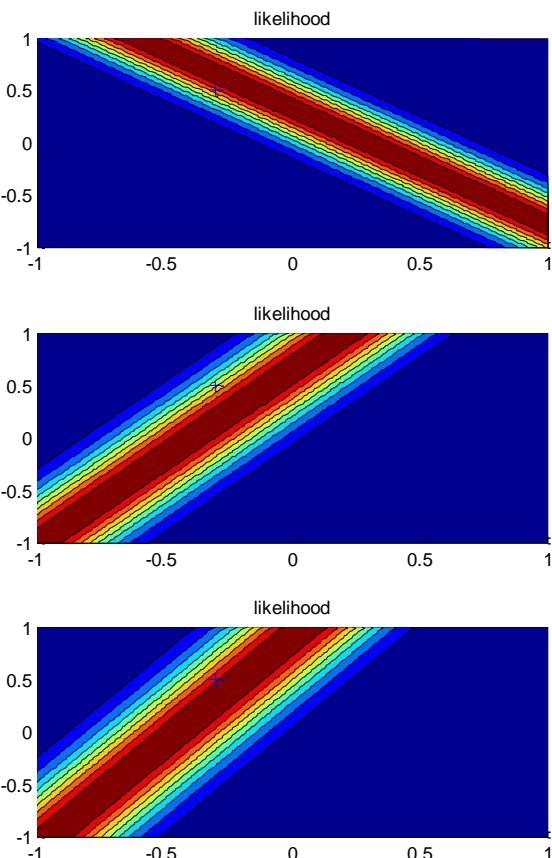


Note that the regression lines after 22 data points have been collected



# Summary of Results

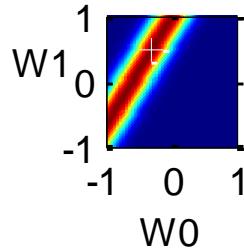
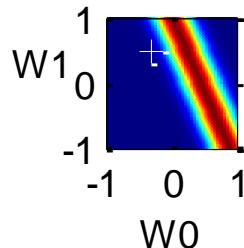
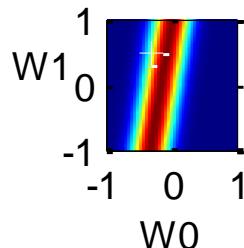
## MatLab Code



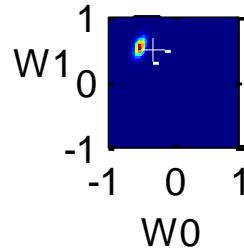
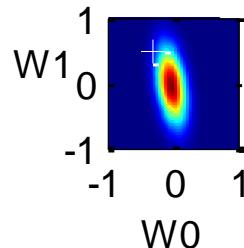
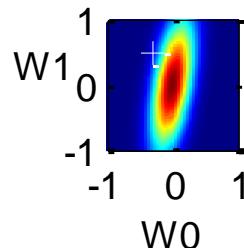
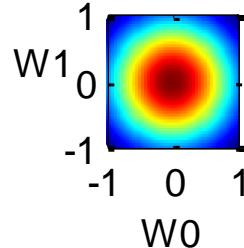
# Summary of Results

likelihood

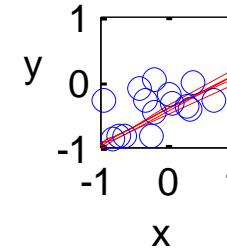
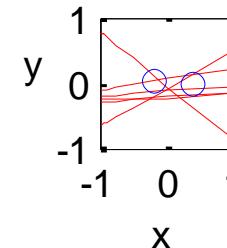
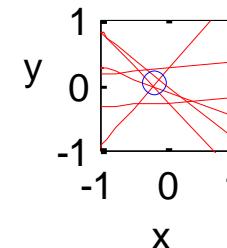
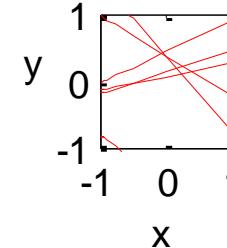
Run [bayesLinRegDemo2d](#)  
from PMTK3



prior/posterior



data space



After 20 data points

# Predictive Distribution

---

- In a full Bayesian treatment, we want to compute the predictive distribution, i.e. given the training data  $\mathbf{x}$  and  $\mathbf{t}$  and a new test point  $x$ , we want the distribution:

$$p(t | \mathbf{x}, \mathbf{x}, \mathbf{t}) = \int p(t | \mathbf{x}, \mathbf{w}) p(\mathbf{w} | \mathbf{x}, \mathbf{t}) d\mathbf{w}, \text{ where}$$

$$p(t | \mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t | y(\mathbf{x}, \mathbf{w}), \beta^{-1}) = \mathcal{N}(t | \boldsymbol{\phi}(\mathbf{x})^T \mathbf{w}, \beta^{-1}) \text{ and}$$

$$p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \alpha, \beta) \propto \mathcal{N}\left(\mathbf{w} / \beta \mathbf{S}_N \sum_{n=1}^N t_n \boldsymbol{\phi}(x_n), \mathbf{S}_N^{-1}\right), \quad \mathbf{S}_N^{-1} = \alpha \mathbf{I} + \sum_{n=1}^N \beta \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T$$

- To integrate  $\mathbf{w}$  out as shown on the predictive distribution expression, we use a fundamental result for linear Gaussian models.

# Appendix: Linear Gaussian Models

$$p(\boldsymbol{x}) = \mathcal{N}(\boldsymbol{x} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$

$$p(\boldsymbol{y} | \boldsymbol{x}) = \mathcal{N}(\boldsymbol{y} | \boldsymbol{A}\boldsymbol{x} + \boldsymbol{b}, \boldsymbol{L}^{-1})$$

- For the above linear Gaussian model, the following very useful results hold:

$$p(\boldsymbol{y}) = \mathcal{N}(\boldsymbol{y} | \boldsymbol{A}\boldsymbol{\mu} + \boldsymbol{b}, \boldsymbol{L}^{-1} + \boldsymbol{A}\boldsymbol{\Lambda}^{-1}\boldsymbol{A}^T)$$

$$p(\boldsymbol{x} | \boldsymbol{y}) = \mathcal{N}(\boldsymbol{x} | (\boldsymbol{\Lambda} + \boldsymbol{A}^T \boldsymbol{L} \boldsymbol{A})^{-1} (\boldsymbol{\Lambda}\boldsymbol{\mu} + \boldsymbol{A}^T \boldsymbol{L}(\boldsymbol{y} - \boldsymbol{b})), (\boldsymbol{\Lambda} + \boldsymbol{A}^T \boldsymbol{L} \boldsymbol{A})^{-1})$$

# Predictive Distribution

$$p(x) = \mathcal{N}(x | \mu, \Lambda^{-1})$$

$$p(y | x) = \mathcal{N}(y | Ax + b, L^{-1})$$

$$p(w | x, t, \alpha, \beta) = \mathcal{N}\left(w | \beta S_N \sum_{n=1}^N t_n \phi(x_n), S_N\right)$$

$$p(t | x, w, \beta) = \mathcal{N}(t | y(x, w), \beta^{-1})$$

➤ Thus for our problem:

$$x \leftarrow w, \mu = \beta S_N \sum_{n=1}^N t_n \phi(x_n), \Lambda^{-1} = S_N$$

$$y \leftarrow t, A = \phi(x)^T, b = 0, L^{-1} = \beta^{-1}$$

➤ The predictive distribution now takes the form:

$$p(y) = \mathcal{N}(y | A\mu + b, L^{-1} + A\Lambda^{-1}A^T) \Rightarrow$$

$$p(t) = \mathcal{N}\left(t | \phi(x)^T \beta S_N \sum_{n=1}^N t_n \phi(x_n), \beta^{-1} + \phi(x)^T S_N \phi(x)\right)$$

# Predictive Distribution

- In a full Bayesian treatment, we want to compute the predictive distribution, i.e. given the training data  $\mathbf{x}$  and  $\mathbf{t}$  and a new test point  $x$ , we want the distribution:

$$p(t | \mathbf{x}, \mathbf{x}, \mathbf{t}) = \int p(t | \mathbf{x}, \mathbf{w}) p(\mathbf{w} | \mathbf{x}, \mathbf{t}) d\mathbf{w}, \quad p(t | \mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t | y(x, \mathbf{w}), \beta^{-1}) \Rightarrow$$

$$p(t | \mathbf{x}, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t | m(x), s^2(x))$$

where the mean and variance were shown in the earlier slide to be:

$$m(x) = \beta \boldsymbol{\phi}(x)^T \mathbf{S}_N \sum_{n=1}^N \boldsymbol{\phi}(x_n) t_n$$

$$s^2(x) = \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S}_N \boldsymbol{\phi}(x)$$

$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \sum_{n=1}^N \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T$$

# Predictive Distribution

- The notation used here is as follows:

$$p(t | x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t | m(x), s^2(x))$$

$$m(x) = \beta \phi(x)^T S_N \sum_{n=1}^N \phi(x_n) t_n$$
$$s^2(x) = \beta^{-1} + \phi(x)^T S_N \phi(x)$$

$$S_N^{-1} = \alpha I + \beta \sum_{n=1}^N \phi(x_n) \phi(x_n)^T$$

Predictive mean and variance are functions of  $x$ .

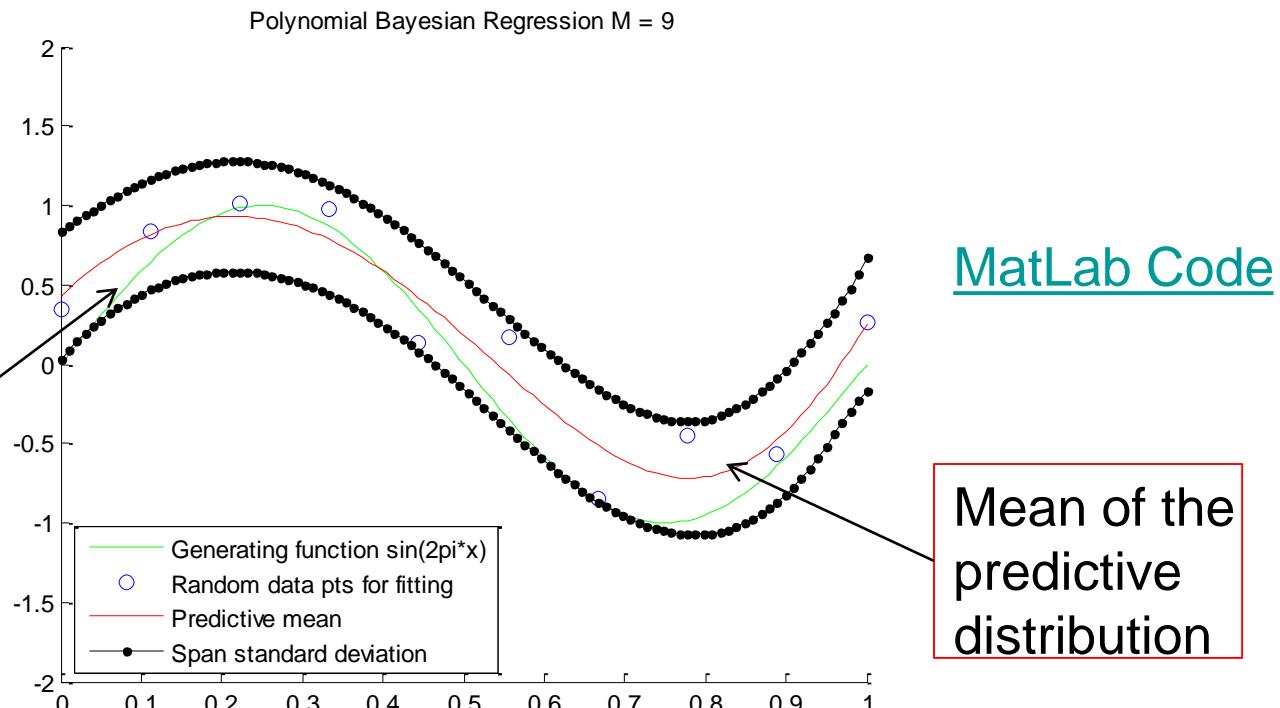
Eg. for polynomial regression:  $\phi(x_n) = \begin{Bmatrix} 1 \\ x_n \\ x_n^2 \\ \vdots \\ x_n^{M-1} \end{Bmatrix}$ ,

$$\phi(x)^T = \{1 \quad x \quad x^2 \quad \dots \quad x^{M-1}\}, I = \text{unit matrix } M \times M$$



# Predictive Distribution

- The predictive distribution using an  $M = 9$  polynomial, with fixed parameters  $\alpha = 5 \times 10^{-3}$  and  $\beta = 11.1$  (corresponding to a known noise variance).
- The red curve denotes the mean of the predictive distribution and the red region corresponds to  $\pm$  standard deviation around the mean.



# Predictive Distribution

- In a full Bayesian treatment, we want to compute the **predictive distribution**, i.e. given the training data  $\mathbf{x}$  and  $\mathbf{t}$  and a new test point  $x$ , we want the distribution:

$$p(t | \mathbf{x}, \mathbf{x}, \mathbf{t}) = \int p(t | \mathbf{x}, \mathbf{w}) p(\mathbf{w} | \mathbf{x}, \mathbf{t}) d\mathbf{w}, \quad p(t | \mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t | y(\mathbf{x}, \mathbf{w}), \beta^{-1}) \Rightarrow$$

$$p(t | \mathbf{x}, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t | m(x), \sigma_N^2(x))$$

where the mean and variance are given by

$$m(x) = \beta \boldsymbol{\phi}(x)^T \mathbf{S}_N \sum_{n=1}^N \boldsymbol{\phi}(x_n) t_n = \boldsymbol{\phi}(x)^T \mathbf{m}_N, \quad \mathbf{m}_N = \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t}$$

$$\sigma_N^2(x) = \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S}_N \boldsymbol{\phi}(x) \text{ (*uncertainty in the data + uncertainty in w*)}$$

$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \sum_{n=1}^N \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T = \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

- Note that:  $\sigma_N^2(x) = \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S}_N \boldsymbol{\phi}(x) \xrightarrow{N \rightarrow \infty} \beta^{-1}$  and  $\sigma_{N+1}^2(x) \leq \sigma_N^2(x)$

# Predictive Distribution

□ It is easy to show:

$$\sigma_N^2(x) = \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S}_N \boldsymbol{\phi}(x) \xrightarrow{N \rightarrow \infty} \beta^{-1} \text{ and } \sigma_{N+1}^2(x) \leq \sigma_N^2(x)$$

□ Note:  $\mathbf{S}_{N+1}^{-1} = \alpha \mathbf{I} + \beta \sum_{n=1}^{N+1} \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T = \mathbf{S}_N^{-1} + \beta \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T$

and the identity:

$$(\mathbf{M} + \mathbf{v} \mathbf{v}^T)^{-1} = \mathbf{M}^{-1} - \frac{(\mathbf{M}^{-1} \mathbf{v})(\mathbf{v}^T \mathbf{M}^{-1})}{1 + \mathbf{v}^T \mathbf{M}^{-1} \mathbf{v}}$$

□ Using these results, we can write:

$$\begin{aligned} \sigma_{N+1}^2(x) &= \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S}_{N+1} \boldsymbol{\phi}(x) = \beta^{-1} + \boldsymbol{\phi}(x)^T \left( \mathbf{S}_N^{-1} + \beta \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T \right)^{-1} \boldsymbol{\phi}(x) = \beta^{-1} + \\ &\boldsymbol{\phi}(x)^T \left[ \mathbf{S}_N - \frac{\beta (\mathbf{S}_N \boldsymbol{\phi}(x_n)) (\boldsymbol{\phi}(x_n)^T \mathbf{S}_N)}{1 + \beta \boldsymbol{\phi}(x_n)^T \mathbf{S}_N \boldsymbol{\phi}(x_n)} \right] \boldsymbol{\phi}(x) = \sigma_N^2(x) - \frac{\beta (\boldsymbol{\phi}(x_n)^T \mathbf{S}_N \boldsymbol{\phi}(x))^2}{1 + \beta \boldsymbol{\phi}(x_n)^T \mathbf{S}_N \boldsymbol{\phi}(x_n)} \leq \sigma_N^2(x) \end{aligned}$$

# Predictive Distribution: Summary

- The notation used here is as follows:

$$p(t | x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t | m(x), \sigma_N^2(x))$$

$$\begin{aligned}m(x) &= \beta \boldsymbol{\phi}(x)^T \mathbf{S}_N \sum_{n=1}^N \boldsymbol{\phi}(x_n) t_n \\ \sigma_N^2(x) &= \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S}_N \boldsymbol{\phi}(x) \\ \mathbf{S}_N^{-1} &= \alpha \mathbf{I} + \beta \sum_{n=1}^N \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T\end{aligned}$$

Note:  
Predictive mean and variance are functions of  $x$ .

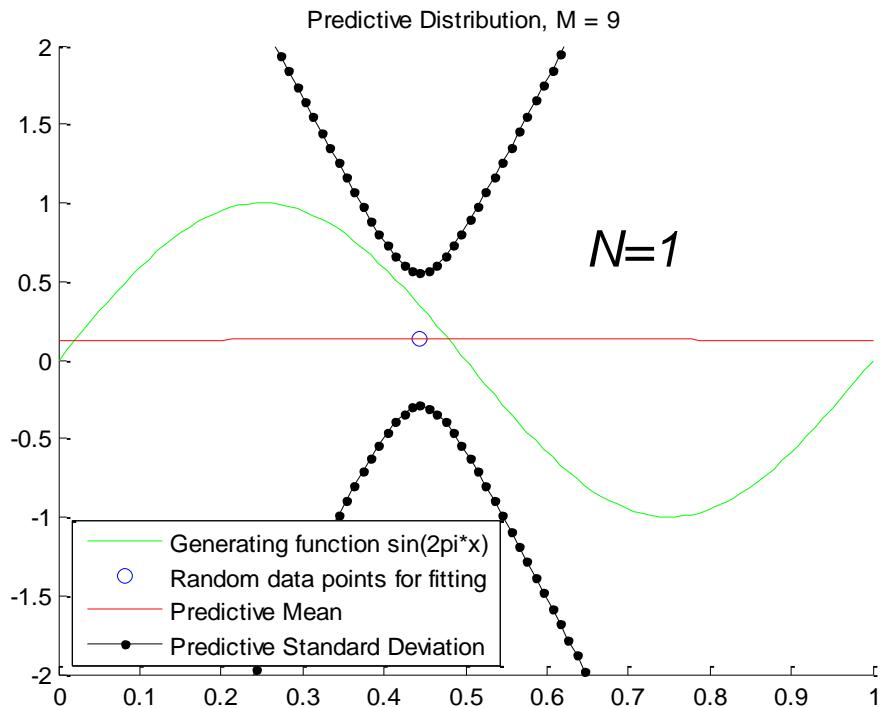
$$\boldsymbol{\phi}(x_n) = \begin{pmatrix} \phi_0(x_n) \\ \phi_1(x_n) \\ \phi_2(x_n) \\ \vdots \\ \phi_{M-1}(x_n) \end{pmatrix}, \boldsymbol{\phi}(x)^T = \{\phi_0(x_n) \quad \phi_1(x_n) \quad \phi_2(x_n) \quad \dots \quad \phi_{M-1}(x_n)\}, \mathbf{I} = \text{unit matrix } M \times M$$

For Polynomial regression:

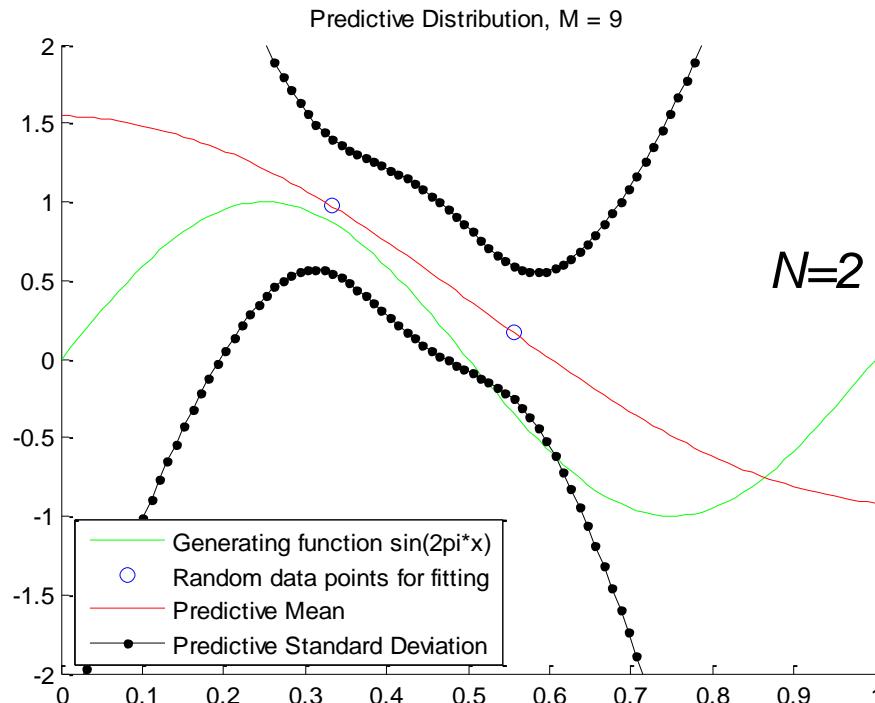
$$\boldsymbol{\phi}(x_n) = \begin{pmatrix} 1 \\ x_n \\ x_n^2 \\ \vdots \\ x_n^{M-1} \end{pmatrix}, \boldsymbol{\phi}(x)^T = \{1 \quad x \quad x^2 \quad \dots \quad x^{M-1}\}$$



# Pointwise Uncertainty in the Predictions



$M=9$  Gaussians, 10 parameters  
Scale of Gaussians adjusted  
with data  
 $\alpha = 5 \times 10^{-3}$   
 $\beta = 11.1$   
Using  $N = 1, 2, 4, 10$   
Data are given [here](#)

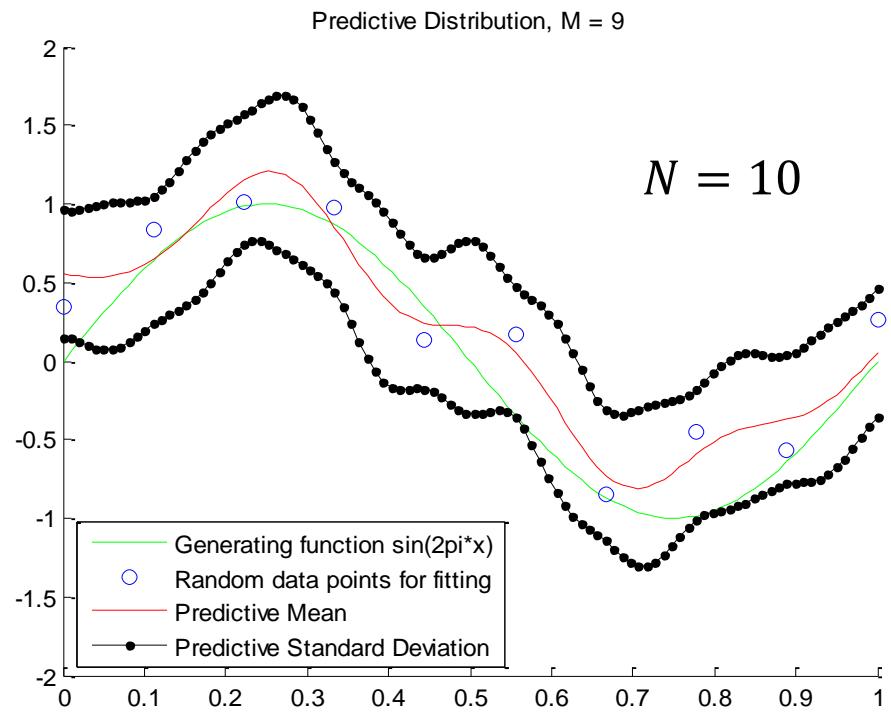
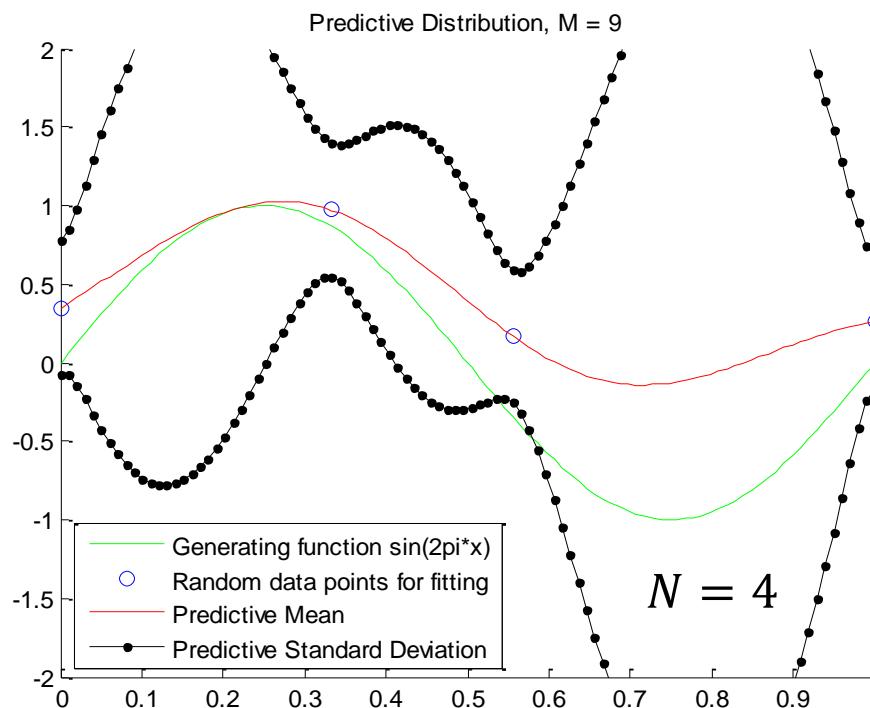


- The predictive uncertainty is smaller near the data.
- The level of uncertainty decreases with  $N$

[MatLab code](#)



# Pointwise Uncertainty in the Predictions

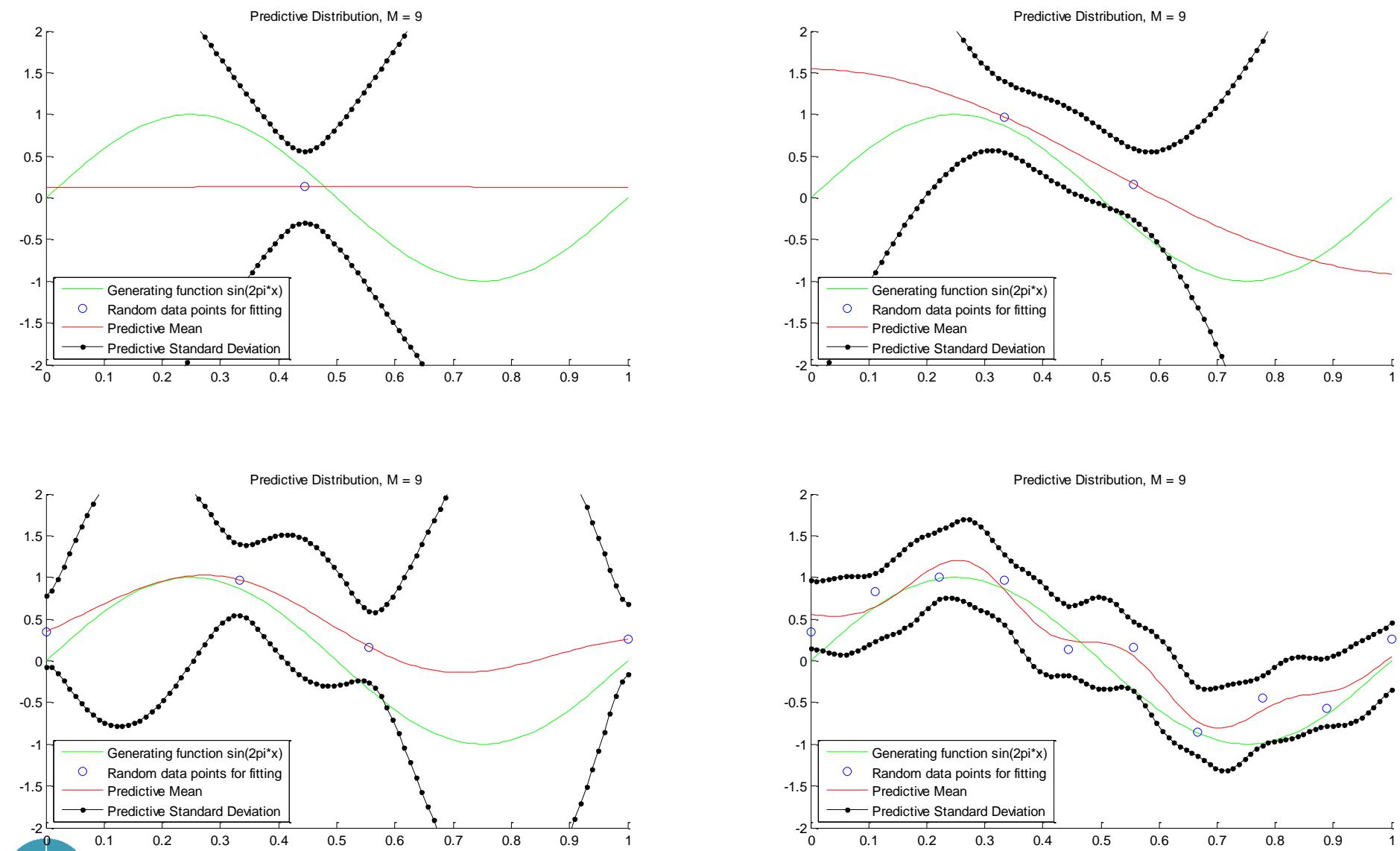


[MatLab code](#)



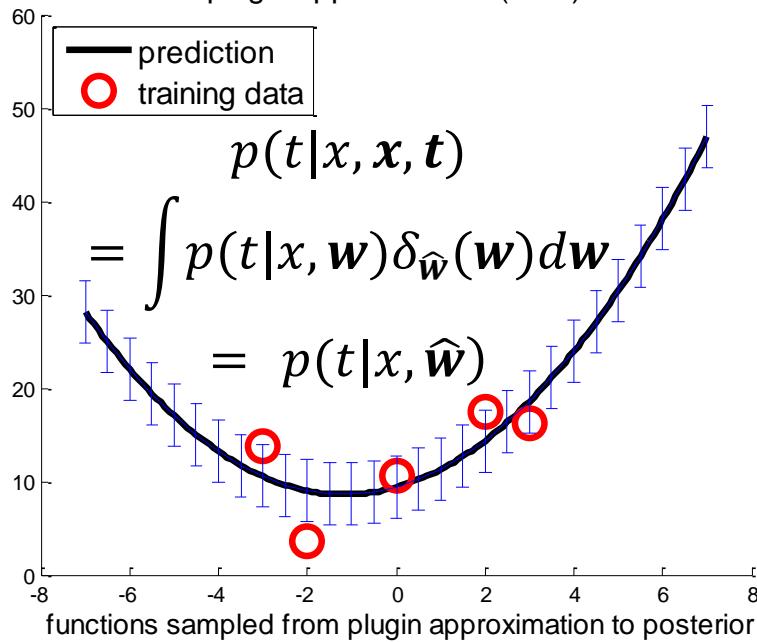
# Summary of Results

## MatLab code

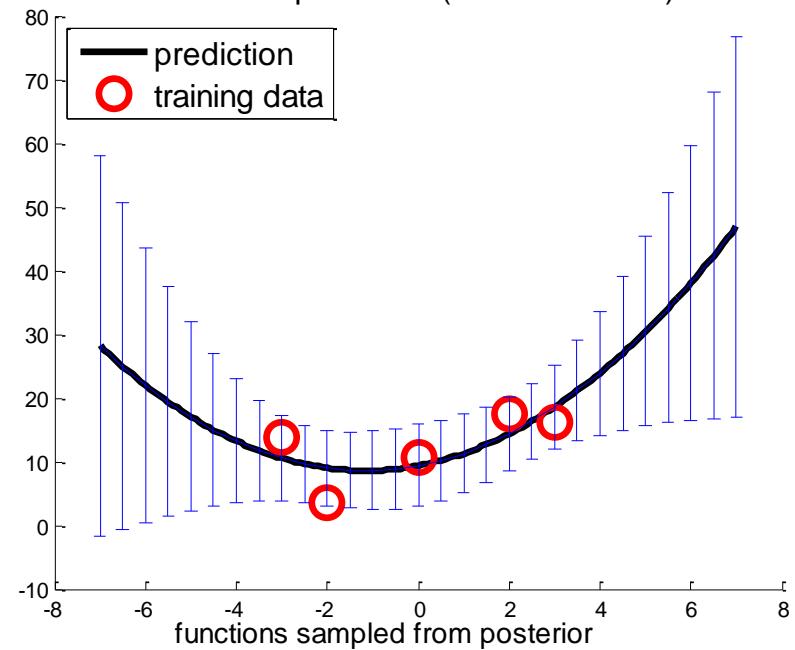


# Plugin Approximation

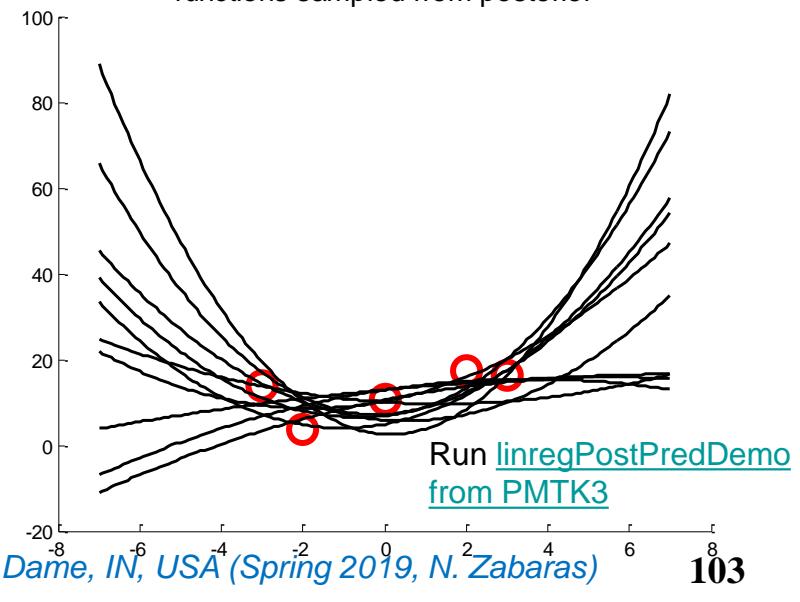
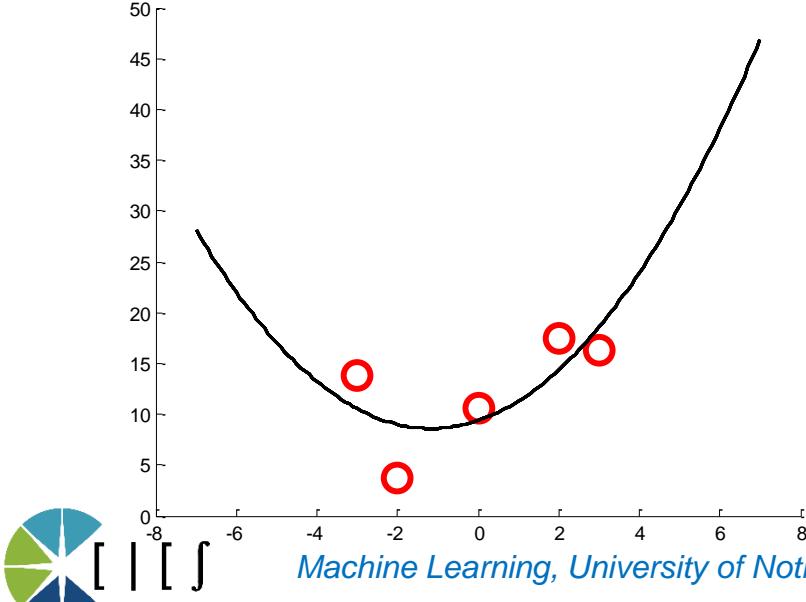
plugin approximation (MLE)



Posterior predictive (known variance)



functions sampled from plugin approximation to posterior



# Predictive Distribution

- We are not interested in  $w$  itself but in making predictions of  $t$  for new values of  $x$ . This requires the predictive distribution

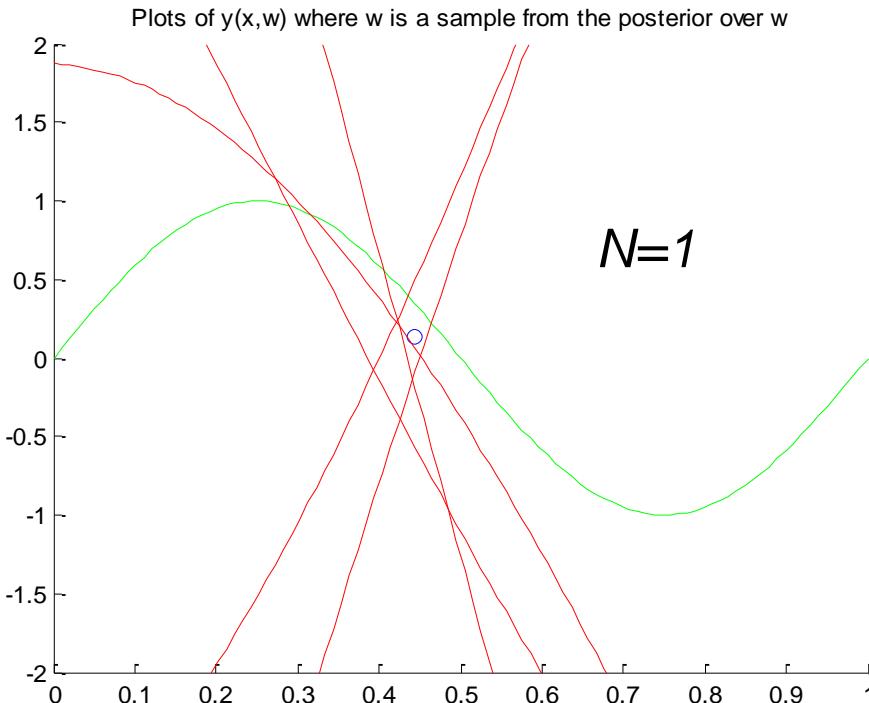
$$p(t | x, \mathbf{x}, \mathbf{t}) = \int p(t | x, w, \beta) p(w | \mathbf{x}, \mathbf{t}, \alpha, \beta) dw = \mathcal{N}\left(t | \mathbf{m}_N^T \phi(x), \sigma_N^2(x)\right)$$

where  $\sigma_N^2(x) = \frac{1}{\beta} + \phi(x)^T S_N \phi(x), \quad S_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi$

- The 1<sup>st</sup> term represents the noise on the data whereas the 2<sup>nd</sup> term reflects the uncertainty associated with  $w$ .
- Because the noise process and the distribution of  $w$  are independent Gaussians, their variances are additive.
- The error bars get larger as we move away from the training points. By contrast, in the plug-in approximation, the error bars are of constant size.
- As additional data points are observed, the posterior distribution becomes narrower.



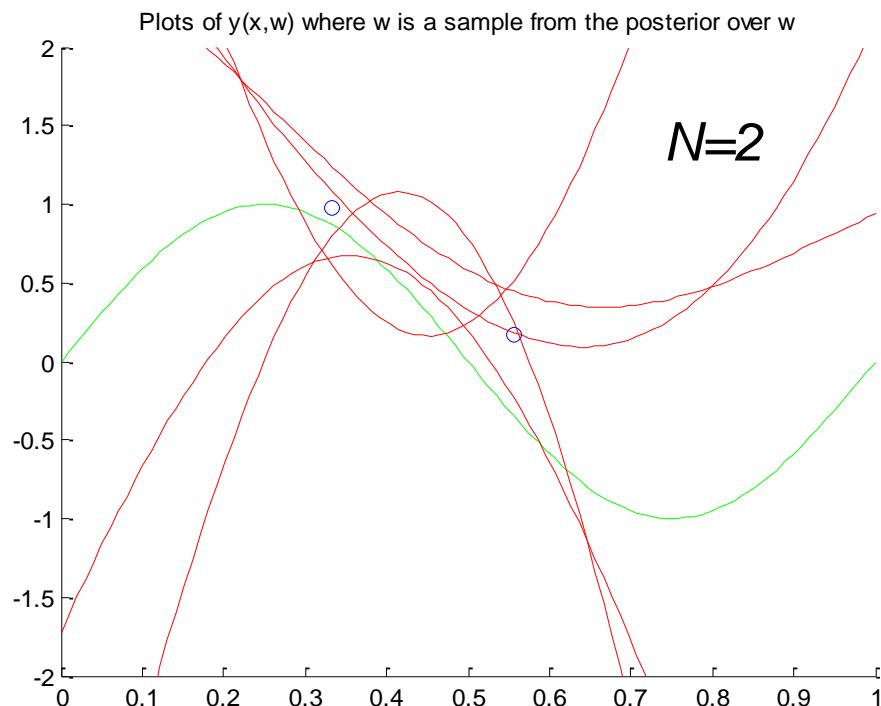
# Covariance Between the Predictions



Same data and basis functions  
as in the earlier example

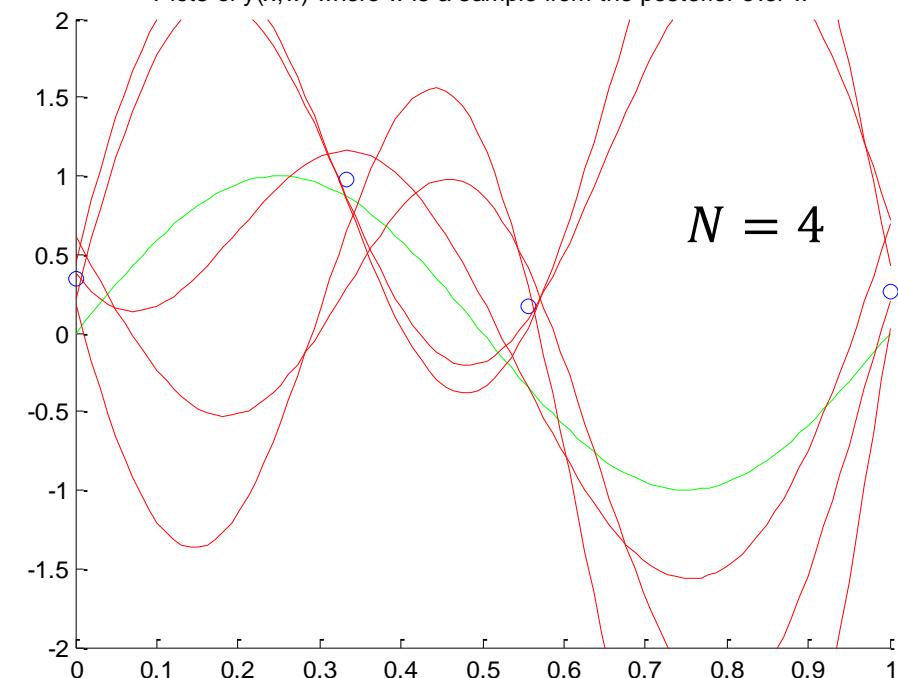
MatLab Code

- Draw samples from the posterior of  $w$  and then plot  $y(x, w)$ . We use the same data as the earlier example.
- We are visualizing the joint uncertainty in the posterior distribution between the  $y$  values at two or more  $x$  values.



# Covariance Between the Predictions

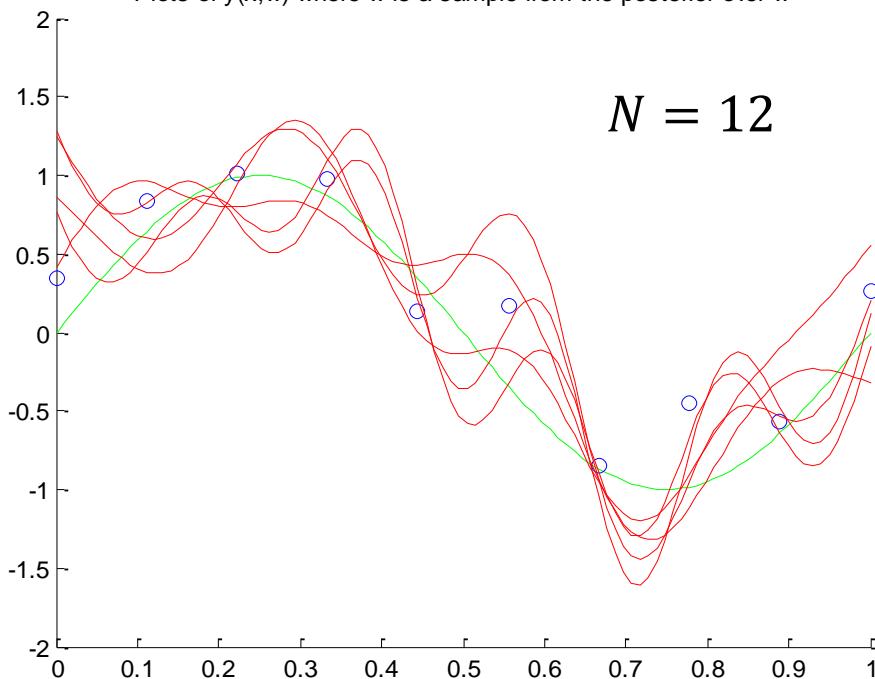
Plots of  $y(x, w)$  where  $w$  is a sample from the posterior over  $w$



$N = 4$

- Draw samples from the posterior of  $w$  and then plot  $y(x, w)$
- We are visualizing the joint uncertainty in the posterior distribution between the  $y$  values at two or more  $x$  values.

Plots of  $y(x, w)$  where  $w$  is a sample from the posterior over  $w$



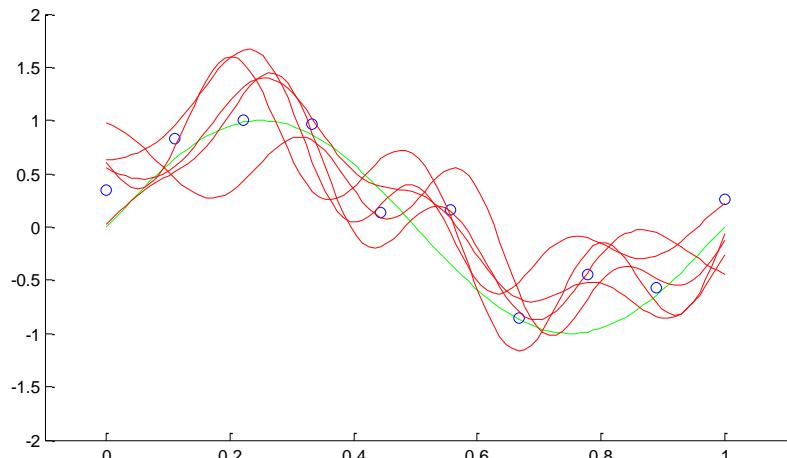
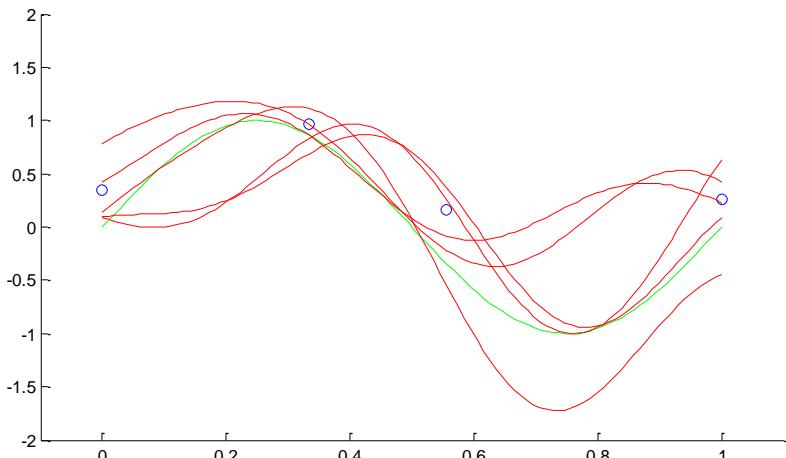
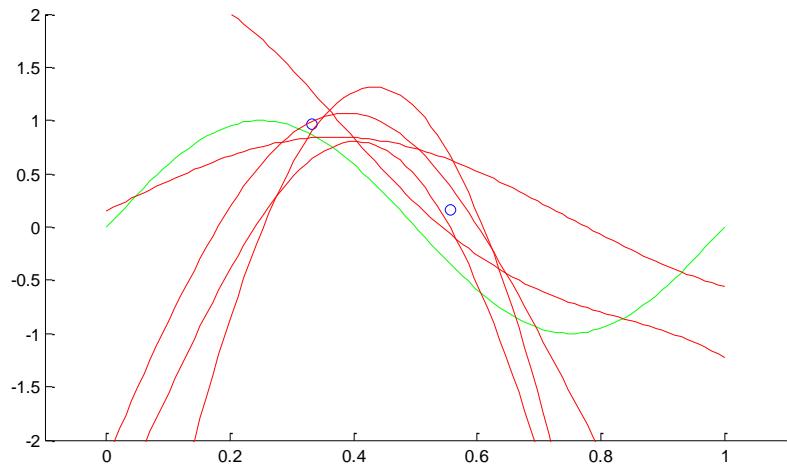
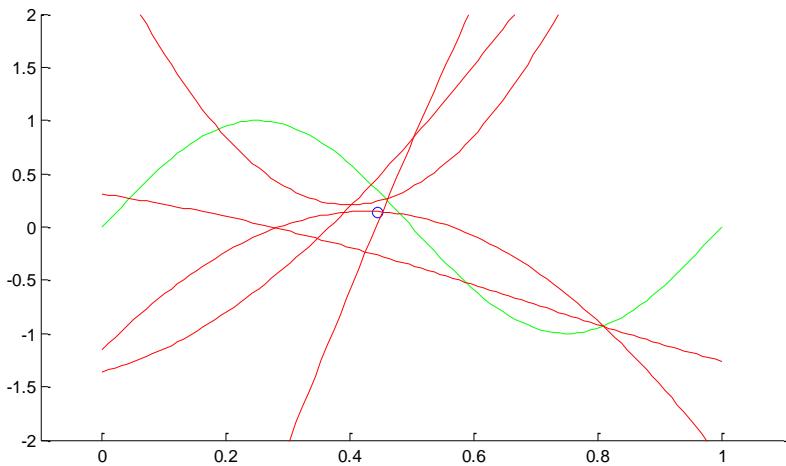
$N = 12$

MatLab Code



# Summary of Results

## MatLab Code



# Gaussian Basis vs. Gaussian Process

- If we use localized basis functions such as Gaussians, then in regions away from the basis function support, the contribution from the second term in the predictive variance will go to zero, leaving only the noise contribution  $\beta^{-1}$ .

$$\sigma_N^2(x) = \beta^{-1} + \phi(x)^T S_N \phi(x) \xrightarrow[\text{support of } \phi(x)]{\text{away from the}} \beta^{-1}$$

- The model becomes very confident in its predictions when extrapolating outside the region occupied by the basis functions. This is an undesirable behavior.
- This problem can be avoided by adopting an alternative Bayesian approach to regression ([Gaussian processes](#)).



# Equivalent Kernel

- The predictive mean in our regression model can be written as:

$$y(\mathbf{x}, \mathbf{m}_N) = \boldsymbol{\phi}(\mathbf{x})^T \mathbf{m}_N = \beta \boldsymbol{\phi}(\mathbf{x})^T S_N \boldsymbol{\Phi}^T \mathbf{t} = \sum_{n=1}^N \underbrace{\beta \boldsymbol{\phi}(\mathbf{x})^T S_N \boldsymbol{\phi}(\mathbf{x}_n)}_{\text{Equivalent kernel } k(\mathbf{x}, \mathbf{x}_n)} t_n \Rightarrow$$

$$y(\mathbf{x}, \mathbf{m}_N) = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n$$

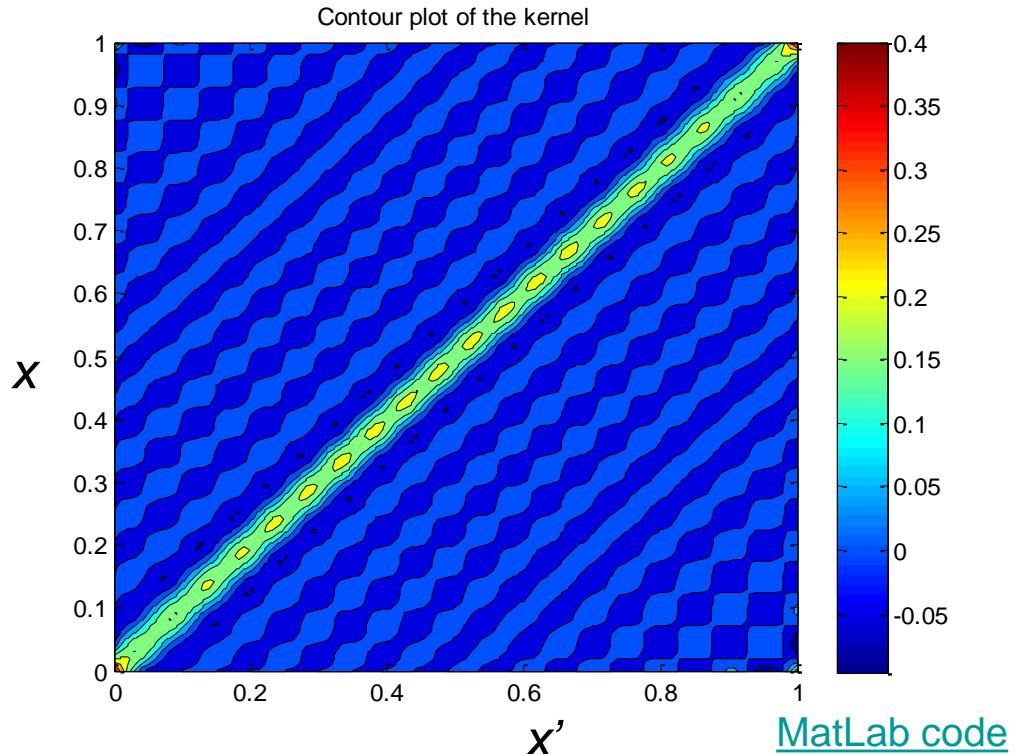
$$k(\mathbf{x}, \mathbf{x}_n) = \beta \boldsymbol{\phi}(\mathbf{x})^T S_N \boldsymbol{\phi}(\mathbf{x}_n)$$

The kernel is shown here as a plot of  $\mathbf{x}$  vs.  $\mathbf{x}'$ .

20 samples  $\mathbf{x}_n$  equally spaced in  $[0,1]$

Gaussian kernels ( $s = 0.05$ ),  $\alpha = 5 \times 10^{-3}$ ,  $\beta = 11.1$ .

$$\phi_j(x) = \exp\left(-\frac{(x - x_j)^2}{2s^2}\right)$$



# Equivalent Kernel

- The predictive mean can be written as follows:

$$y(\mathbf{x}, \mathbf{m}_N) = \boldsymbol{\phi}(\mathbf{x})^T \mathbf{m}_N = \beta \boldsymbol{\phi}(\mathbf{x})^T S_N \Phi^T \mathbf{t} = \sum_{n=1}^N \underbrace{\beta \boldsymbol{\phi}(\mathbf{x})^T S_N \boldsymbol{\phi}(\mathbf{x}_n)}_{\text{Equivalent kernel } k(\mathbf{x}, \mathbf{x}_n)} t_n \Rightarrow$$

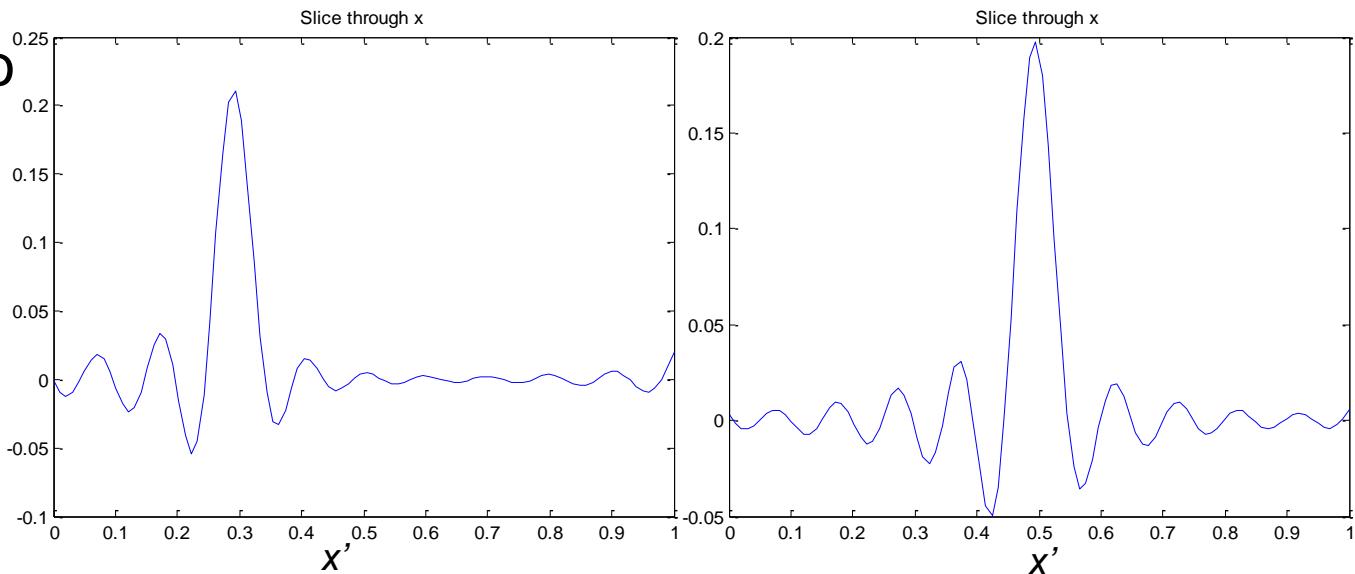
$$y(\mathbf{x}, \mathbf{m}_N) = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n$$

$$k(\mathbf{x}, \mathbf{x}_n) = \beta \boldsymbol{\phi}(\mathbf{x})^T S_N \boldsymbol{\phi}(\mathbf{x}_n)$$

Two slices  
corresponding to  
 $x = 0.3, x = 0.5.$

20 Gaussians  
( $s = 0.05$ )

Plots for 100  
points  $\mathbf{x}$   
uniformly  
in  $(0,1)$



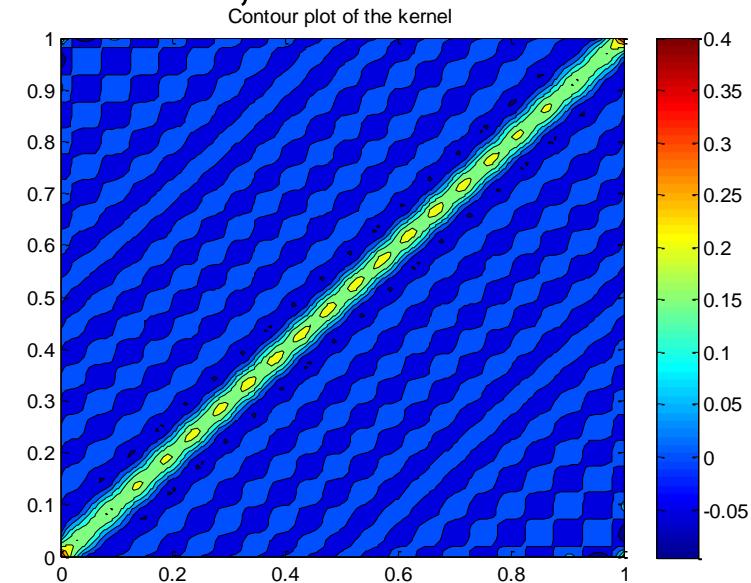
[MatLab code](#)

# Equivalent Kernel

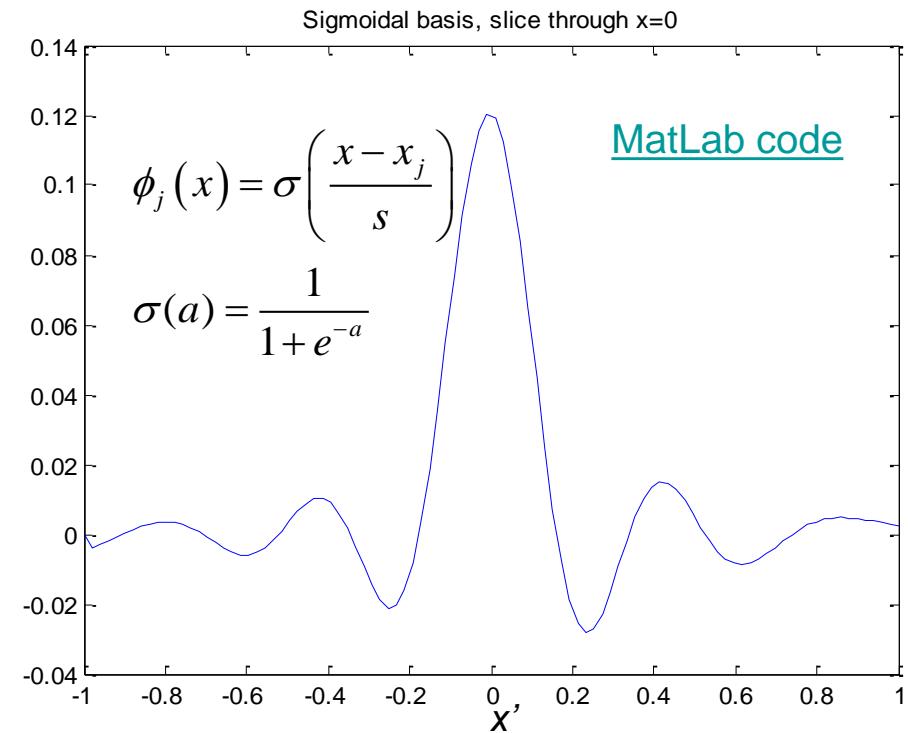
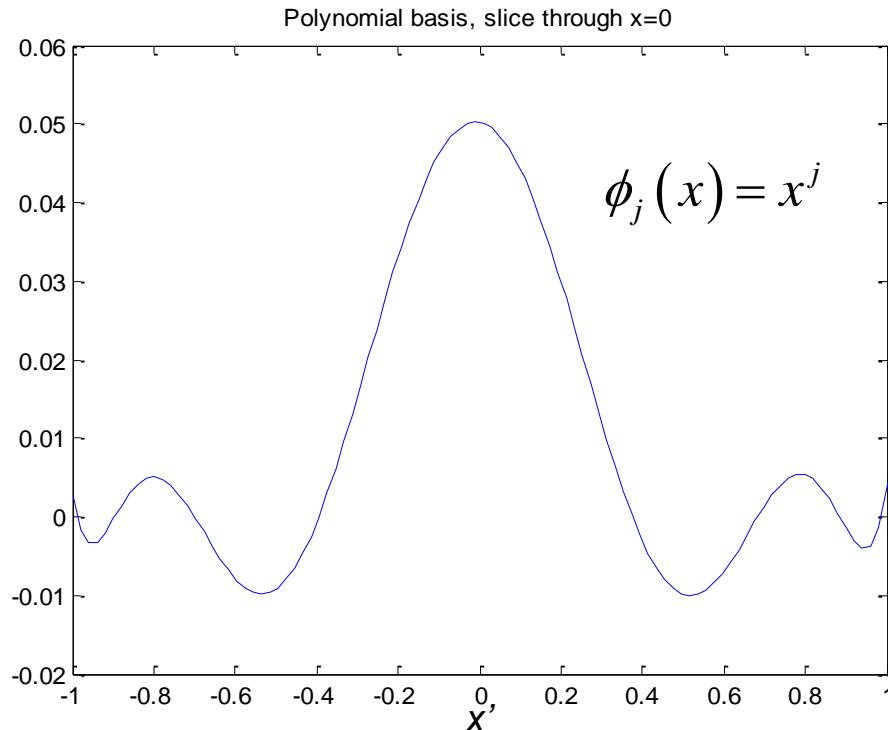
$$y(\mathbf{x}, \mathbf{m}_N) = \sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) t_n$$

$$k(\mathbf{x}, \mathbf{x}_n) = \beta \phi(\mathbf{x})^T S_N \phi(\mathbf{x}_n)$$

- Note that we make predictions by taking linear combinations of the training set target values (*linear smoothers*).
- The equivalent kernel depends on the input values  $\mathbf{x}_n$  since these appear in  $S_N$ .
- See from the figure that its localized around  $\mathbf{x}$ , and so the mean of the predictive distribution  $y(\mathbf{x}, \mathbf{m}_N)$ , is obtained by forming a weighted combination of  $t_n$  in which data points close to  $\mathbf{x}$  are given higher weight than points (evidence) further removed from  $\mathbf{x}$ .



# Equivalent Kernel



- ❑ Examples of equivalent kernels  $k(x, x')$  for  $x = 0$  plotted as a function of  $x'$ . 10 basis functions were used in each case. The sigmoidal basis is centered at 10 equally spaced  $x_n$  in  $[-1,1]$ .
- ❑ These are localized functions of  $x'$  even though the corresponding basis functions are nonlocal.

# Equivalent Kernel

- Consider the covariance between  $y(\mathbf{x})$  and  $y(\mathbf{x}')$ , which is given by (recall that  $p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \alpha, \beta) = \mathcal{N}\left(\mathbf{w} | \mathbf{\mu}_\mathbf{w}, \beta S_N \sum_{n=1}^N t_n \phi(\mathbf{x}_n), S_N\right)$ ):
$$\text{cov}[y(\mathbf{x}), y(\mathbf{x}')] = \text{cov}[\phi(\mathbf{x})^T \mathbf{w}, \mathbf{w}^T \phi(\mathbf{x}')] = \phi(\mathbf{x})^T S_N \phi(\mathbf{x}') = \beta^{-1} k(\mathbf{x}, \mathbf{x}')$$
- From [the earlier discussion](#) on the equivalent kernel, we see that [the predictive mean at nearby points will be highly correlated](#), whereas for more distant points the correlation will be smaller.
- By drawing samples from the posterior distribution over  $\mathbf{w}$ , and plotting the corresponding model functions  $y(\mathbf{x}, \mathbf{w})$ , [we are visualizing](#) the joint uncertainty in the posterior distribution between the  $y$  values at two (or more)  $\mathbf{x}$  values, as governed by the equivalent kernel.
- This is contrary to the [Figure here](#) where we visualized pointwise uncertainty in the predictions.



# Equivalent Kernel

---

- We can avoid the use of basis functions and define the kernel function directly, leading e.g. to *Gaussian Processes* (more on Kernel methods later on in this course).
- Note that

$$\sum_{n=1}^N k(\mathbf{x}, \mathbf{x}_n) = 1$$

for all values of  $\mathbf{x}$ . However, the equivalent kernel may be negative for some values of  $\mathbf{x}$ .

- Like all kernel functions, the equivalent kernel can be expressed as an inner product:

$$k(\mathbf{x}, \mathbf{z}) = \beta \phi(\mathbf{x})^T S_N \phi(\mathbf{z}) \Rightarrow$$

$$k(\mathbf{x}, \mathbf{z}) = \psi(\mathbf{x})^T \psi(\mathbf{z})$$

$$\psi(\mathbf{x}) = \beta^{1/2} S_N^{1/2} \phi(\mathbf{x})$$



# Computing the Bias Parameter

- If we make the bias parameter  $w_0$  explicit, then the error function becomes

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left( t_n - w_0 - \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}_n) \right)^2$$

- Setting the derivative with respect to  $w_0$  equal to zero, and solving for  $w_0$ , we obtain

$$w_0 = \frac{1}{N} \sum_{n=1}^N t_n - \frac{1}{N} \sum_{n=1}^N \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}_n) \Rightarrow$$

$$w_0 = \bar{t} - \sum_{j=1}^{M-1} w_j \bar{\phi}_j, \quad \bar{t} = \frac{1}{N} \sum_{n=1}^N t_n, \quad \bar{\phi}_j = \frac{1}{N} \sum_{n=1}^N \phi_j(\mathbf{x}_n)$$

- The bias parameter  $w_0$  compensates for the difference between the averages of the target values and the weighted sum of the averages of the basis function values.

# A Note on Data Centering

- In linear regression, it helps to center the data in a way that does not require us to compute the offset term  $w_0 \equiv \mu$ . Write the likelihood as:

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \mu, \beta) \propto \exp\left(-\frac{\beta}{2} (\mathbf{t} - \mu \mathbf{1}_N - \Phi \mathbf{w})^T (\mathbf{t} - \mu \mathbf{1}_N - \Phi \mathbf{w})\right)$$

- Let us assume that the input data are centered in each dimension such that:

$$\sum_{j=1}^N \phi_i(x_j) = 0 \quad \forall i = 1, \dots, M-1$$
$$\Phi = \begin{pmatrix} \phi^T(x_1) \\ \phi^T(x_2) \\ \vdots \\ \phi^T(x_N) \end{pmatrix} = \begin{pmatrix} \phi_1(x_1) & \phi_2(x_1) & \dots & \phi_{M-1}(x_1) \\ \phi_1(x_2) & \phi_2(x_2) & \dots & \phi_{M-1}(x_2) \\ \vdots & \vdots & \vdots & \vdots \\ \phi_1(x_N) & \phi_2(x_N) & \dots & \phi_{M-1}(x_N) \end{pmatrix}, \quad \begin{aligned} \Phi^T &= (\phi(x_1) \quad \phi(x_2) \quad \dots \quad \phi(x_N)) \\ \phi(x_i) &\equiv (\phi_1(x_i) \quad \phi_2(x_i) \quad \dots \quad \phi_{M-1}(x_i))^T \\ \Phi &= (\varphi_1 \quad \varphi_2 \quad \dots \quad \varphi_{M-1}) \\ \varphi_i &\equiv (\phi_i(x_1) \quad \phi_i(x_2) \quad \dots \quad \phi_i(x_N))^T \end{aligned}$$

- The mean of the output is equally likely to be positive or negative. Let us put an improper prior  $p(\mu) \propto 1$  and integrate  $\mu$  out.



# A Note on Data Centering

- Introducing,  $\bar{t} = \frac{1}{N} \sum_{i=1}^N t_i$ , the marginal likelihood becomes :

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) \propto \int \exp \left( -\frac{\beta}{2} \left( \underbrace{\mathbf{t} - \bar{t}\mathbf{1}_N - \Phi\mathbf{w}}_A - (\mu - \bar{t})\mathbf{1}_N \right)^T \left( \mathbf{t} - \bar{t}\mathbf{1}_N - \Phi\mathbf{w} - (\mu - \bar{t})\mathbf{1}_N \right) \right) d\mu$$

- Completing the square in  $\mu$  gives (use the centering of  $\Phi$ ):

$$\begin{aligned} p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) &\propto \int \exp \left( -\frac{\beta}{2} \left( \begin{array}{c} (\mu - \bar{t})^2 N - 2(\mu - \bar{t}) \\ \underbrace{\mathbf{A}^T \mathbf{1}_N}_{\bar{t}N - \bar{t}N - \underbrace{\mathbf{w}^T \Phi^T \mathbf{1}_N}_{=0}} + \mathbf{A}^T \mathbf{A} \end{array} \right) \right) d\mu \\ &\propto \exp \left( -\frac{\beta}{2} \left( \mathbf{t} - \bar{t}\mathbf{1}_N - \Phi\mathbf{w} \right)^T \left( \mathbf{t} - \bar{t}\mathbf{1}_N - \Phi\mathbf{w} \right) \right) \end{aligned}$$

- Our model is now simplified if instead of  $\mathbf{t}$  we use (centered output)  $\mathbf{t} \leftarrow \mathbf{t} - \bar{t}\mathbf{1}_N$  and the likelihood is simply written as:

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) \propto \exp \left( -\frac{\beta}{2} (\mathbf{t} - \Phi\mathbf{w})^T (\mathbf{t} - \Phi\mathbf{w}) \right)$$

$\hat{\mu} = \bar{t} - \sum_{j=1}^M \bar{\Phi}_j^T \mathbf{w}_j$ ,  $[\bar{\Phi}_1, \dots, \bar{\Phi}_{M-1}]$  is formed

- Recall that the MLE estimate for  $\mu$  is: *by taking the average of each column of  $\Phi$*

# A Note on Data Centering

- As an example, consider a linear regression model of the form

$$\mathbb{E}[y | \mathbf{x}] = w_0 + \mathbf{w}^T \mathbf{x}$$

- In the context e.g. of MLE, we need to minimize

$$\min_{w_0, \mathbf{w}} = \sum_{i=1}^N (t_i - w_0 - \mathbf{w}^T \mathbf{x}_i)^2$$

- Minimization wrt  $w_0$  gives:

$$\sum_{i=1}^N (t_i - w_0 - \mathbf{w}^T \mathbf{x}_i) = 0 \Rightarrow w_0 N = \bar{t} N - N \mathbf{w}^T \bar{\mathbf{x}} \Rightarrow w_0 = \bar{t} - \mathbf{w}^T \bar{\mathbf{x}}$$

where:

$$\bar{\mathbf{x}} \equiv \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_M \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^N x_{i1} / N \\ \sum_{i=1}^N x_{i2} / N \\ \vdots \\ \sum_{i=1}^N x_{iM} / N \end{pmatrix}, \quad \bar{t} \equiv \sum_{i=1}^N t_i / N$$

- Thus:

$$\hat{w}_0 = \bar{t} - \bar{\mathbf{x}}^T \mathbf{w}$$



# A Note on Data Centering

- Substituting the bias term in our objective function gives:

$$\min_{\mathbf{w}} \sum_{i=1}^N \left( t_i - \bar{t} + \mathbf{w}^T \bar{\mathbf{x}} - \mathbf{w}^T \mathbf{x}_i \right)^2 = \min_{\mathbf{w}} \sum_{i=1}^N \left( t_i - \bar{t} + \mathbf{w}^T (\bar{\mathbf{x}} - \mathbf{x}_i) \right)^2$$

- Minimization wrt  $\mathbf{w}$  gives:

$$\sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \hat{\mathbf{w}} = \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{t}_i - \bar{\mathbf{t}})$$

- We thus first compute the MLE of  $\mathbf{w}$  using the centered input and output as follows:

$$\hat{\mathbf{w}} = (\mathbf{X}_c^T \mathbf{X}_c)^{-1} \mathbf{X}_c^T \mathbf{t}_c = \left[ \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \right] \left[ \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{t}_i - \bar{\mathbf{t}}) \right]$$
$$\mathbf{X}_c = \mathbf{X} - \bar{\mathbf{X}} = \mathbf{X} - \mathbf{1}_N \bar{\mathbf{x}}^T = \begin{pmatrix} \mathbf{x}_1^T - \bar{\mathbf{x}}^T \\ \mathbf{x}_2^T - \bar{\mathbf{x}}^T \\ \vdots \\ \mathbf{x}_N^T - \bar{\mathbf{x}}^T \end{pmatrix} = \begin{pmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \dots & x_{1M} - \bar{x}_M \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & \dots & x_{2M} - \bar{x}_M \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} - \bar{x}_1 & x_{N2} - \bar{x}_2 & \dots & x_{NN} - \bar{x}_M \end{pmatrix}, \quad \bar{\mathbf{x}} \equiv \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_M \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^N x_{i1} / N \\ \sum_{i=1}^N x_{i2} / N \\ \vdots \\ \sum_{i=1}^N x_{iM} / N \end{pmatrix}, \quad \mathbf{t}_c = \mathbf{t} - \bar{\mathbf{t}} = \mathbf{t} - \mathbf{1}_N \bar{\mathbf{t}}, \quad \bar{\mathbf{t}} \equiv \sum_{i=1}^N t_i / N$$

- We can then estimate the MLE estimate of  $w_0$  as follows:

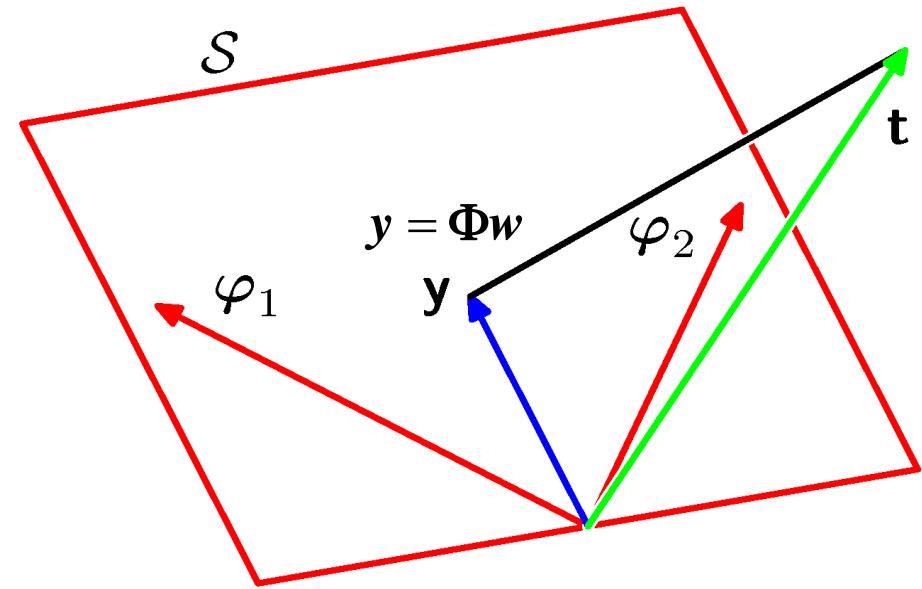
$$\hat{w}_0 = \bar{t} - \bar{\mathbf{x}}^T \mathbf{w}$$



# Geometry of Least Squares

- We look for a geometrical interpretation of the least-squares solution in an  $N$ -dimensional space.  $t$  is a vector in that space with components  $t_1, \dots, t_N$  ( $N > M$ ).
- The least-squares regression function is obtained by finding the orthogonal projection of the data vector  $t$  onto the subspace spanned by the basis functions  $\varphi_j(x)$ .
- Note  $\varphi_j$  is here the  $j$  –th column of  $\Phi$ .

$$\varphi_j \equiv (\phi_j(x_1), \dots, \phi_j(x_N))^T$$



$$\Phi = (\varphi_0 \quad \varphi_1 \quad \dots \quad \varphi_{M-1}) = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \dots & \phi_{M-1}(x_2) \\ \vdots & \vdots & \vdots & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \dots & \phi_{M-1}(x_N) \end{pmatrix}$$

# Geometry of Least Squares

- We are looking for  $w$  such that the projection error

$$t - y = t - \Phi w$$

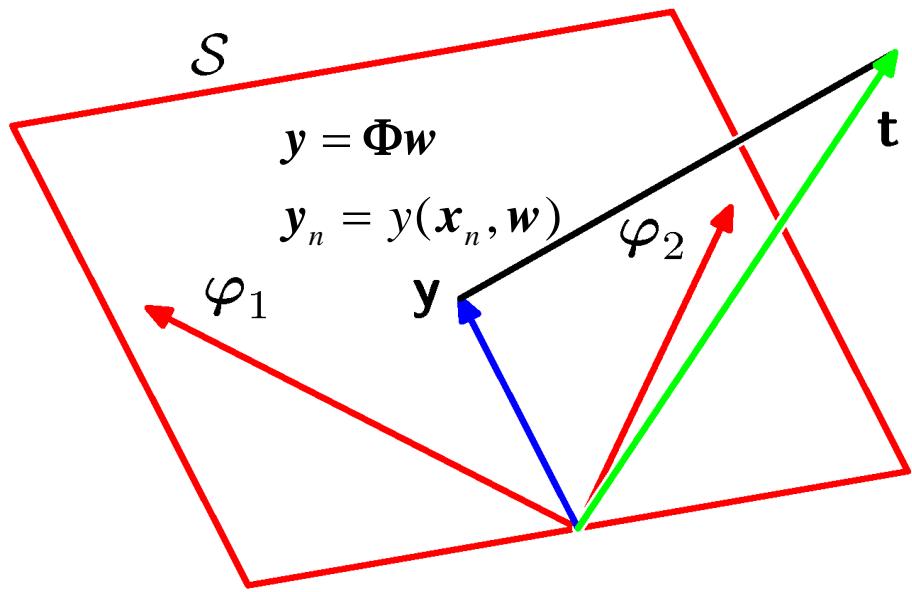
is orthogonal to the basis  $\varphi_j$ , i.e. such that:

$$\Phi^T (t - \Phi w) = 0$$

- These are the normal equations we derived earlier.

$$\Phi^T = \begin{pmatrix} \varphi_0^T \\ \varphi_1^T \\ \vdots \\ \varphi_{M-1}^T \end{pmatrix}$$

$$\varphi_j(x) = (\phi_j(x_1), \dots, \phi_j(x_N))^T$$



$S$ :  $M$  –dimensional subspace spanned by  $\varphi_j(x)$

$$\Phi = [\varphi_0(x) \quad \varphi_1(x) \quad \dots \quad \varphi_{M-1}(x)]$$